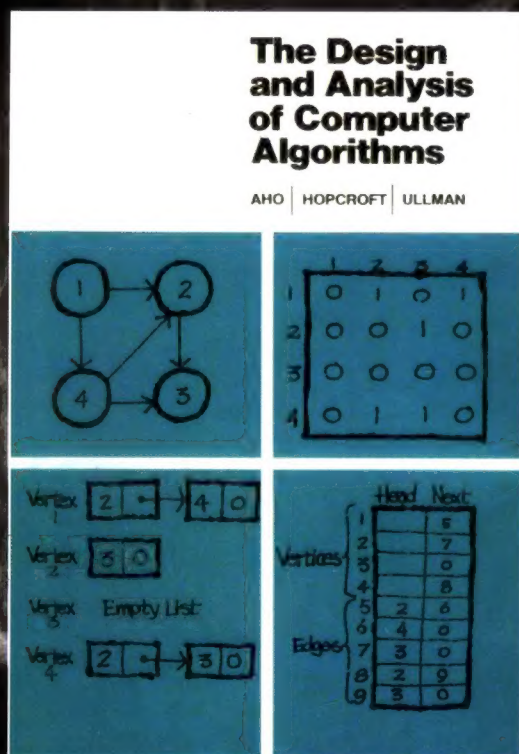


计算机算法的设计与分析

新增经典算法的C/C++实现

(美) Alfred V. Aho John E. Hopcroft Jeffrey D. Ullman 著 黄林鹏 王德俊 张仕 译
哥伦比亚大学 康奈尔大学 斯坦福大学 上海交通大学



The Design and Analysis of Computer Algorithms

计算机算法的设计与分析

本书是著名计算机科学家Alfred V. Aho、John E. Hopcroft和Jeffrey D. Ullman合著的一部经典著作，着重介绍计算机算法设计领域的统一原则和基本概念。书中深入分析一些计算机模型上的算法，以及一些有效算法常用的数据结构和编程技术，为读者提供了有关递归方法、分治方法和动态规划方面的详细实例和实际应用，并致力于更有效算法的设计和开发。同时，对NP完全等问题能否有效求解进行分析，并探索应用启发式算法解决问题的途径。另外，本书还提供了大量富有指导意义的习题。

本书可以作为高等院校计算机专业本科生和研究生算法设计课程的教材，也可以作为计算机算法理论中更高级课程的教材。

本书特点：

- 介绍若干计算模型，包括图灵机模型、随机存取计算机模型以及它们的变体在算法设计和分析中的作用；
- 介绍设计有效算法相关的数据结构，给出递归、分治法、动态规划技术以及相关的算法例子；
- 解释排序算法、集合算法、图算法、模式匹配算法、快速傅里叶变换的技术细节和应用技巧；
- 讨论矩阵乘法、多项式运算以及整数算法的内在机理；
- 探讨一些问题的时间复杂度和空间复杂度，涉及NP完全问题以及一些可以证明不存在有效算法的问题，并将计算难度的概念和向量空间的线性独立性关联在一起；
- 本书算法原采用ALGOL实现，译者为配合国内教学实际，在附录中给出了大多数经典算法的C/C++实现。

作者简介

Alfred V. Aho 博士是哥伦比亚大学计算机科学系主管本科生教学的副主任，IEEE Fellow，美国科学与艺术学院及国家工程院院士，曾获得IEEE的冯·诺伊曼奖。他是《编译原理》（Compiler: Principles, Techniques, and Tools）的第一作者。他目前的研究方向为量子计算、程式设计语言、编译器和算法等。



John E. Hopcroft 博士是康奈尔大学工程学院院长兼计算机科学系教授，IEEE Fellow，美国科学与艺术学院及国家工程院院士，1986年因其在数据结构、算法设计与分析等领域的重要贡献而获得图灵奖。他还是《自动机理论、语言和计算导论》（Introduction to Automata Theory, Languages, and Computation）的第一作者。他目前的研究方向是信息存取。



Jeffrey D. Ullman 博士先后任教于普林斯顿大学和斯坦福大学，现已退休。他是美国国家工程院院士，曾获得1996年的Sigmod贡献奖和2000年的Knuth奖等诸多学术奖项，除本书外，他还与Aho合著了《编译原理》，与Hopcroft合著了《自动机理论、语言和计算导论》，并与其他数据库专家合著了数据库方面的名著，如《数据库系统基础教程》（A First Course in Database Systems）等。



www.PearsonEd.com

投稿热线：(010) 88379604
购书热线：(010) 68995259, 68995264
读者信箱：hzjsj@hzbook.com



华章网站 <http://www.hzbook.com>

网上购书：www.china-pub.com

封面设计：金杨 杨

上架指导：计算机/算法与计算

ISBN 978-7-111-21543-1



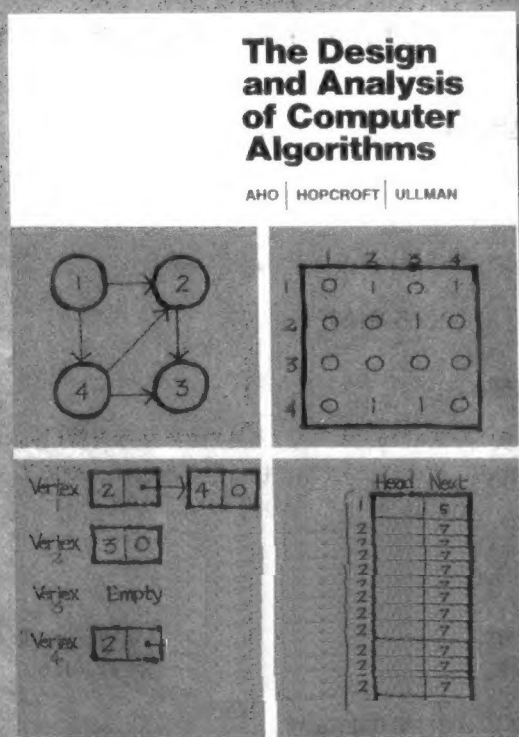
9 787111 215431

ISBN 978-7-111-21543-1
定价：49.00 元

计 算 机 科 学 丛

计算机算法的设计与分析

(美) Alfred V. Aho 约翰·霍普克罗夫特 Jeffrey D. Ullman 著 黄林鹏 王德俊 张仕 译
哥伦比亚大学 康奈尔大学 斯坦福大学 上海交通大学



The Design and Analysis of Computer Algorithms



机械工业出版社
China Machine Press

本书是一部设计与分析领域的经典著作,着重介绍了计算机算法设计领域的基本原则和根本原理。书中深入分析了一些计算机模型上的算法,介绍了一些和设计有效算法有关的数据结构和编程技术,为读者提供了有关递归方法、分治方法和动态规划方面的详细实例和实际应用,并致力于更有效算法的设计和开发。同时,对 NP 完全等问题能否有效求解进行了分析,并探索了应用启发式算法解决问题的途径。另外,本书还提供了大量富有指导意义的习题。

本书可以作为高等院校计算机算法设计与分析课程的本科生或研究生教材,也可以作为计算机理论研究人员、计算机算法设计人员的参考书。

Simplified Chinese edition copyright © 2007 by Pearson Education Asia Limited and China Machine Press.

Original English language title: *The Design and Analysis of Computer Algorithms* (ISBN: 0-201-00029-6)
Alfred V. Aho, John E. Hopcroft, Jeffrey D. Ullman, Copyright © 1974.

All rights reserved.

Published by arrangement with the original publisher, Pearson Education, Inc., publishing as ADDISON WESLEY.

本书封面贴有 Pearson Education (培生教育出版集团) 激光防伪标签,无标签者不得销售。

版权所有,侵权必究。

本书法律顾问 北京市展达律师事务所

本书版权登记号: 图字: 01-2006-3157

图书在版编目 (CIP) 数据

计算机算法的设计与分析/(美)阿霍(Aho, A. V.), (美)霍普克劳夫特(Hopcroft, J. E.), (美)乌尔曼(Ullman, J. D.)著;黄林鹏,王德俊,张仕译.-北京:机械工业出版社,2007.7

(计算机科学丛书)

书名原文: *The Design and Analysis of Computer Algorithms*

ISBN 978-7-111-21543-1

I. 计… II. ①阿… ②霍… ③乌… ④黄… ⑤王… ⑥张… III. ①电子计算机-算法设计-高等学校-教材 ②电子计算机-算法分析-高等学校-教材 IV. TP301.6

中国版本图书馆 CIP 数据核字(2007)第 074676 号

机械工业出版社(北京市西城区百万庄大街 22 号 邮政编码 100037)

责任编辑:王 玉

三河市明辉印装有限公司印刷·新华书店北京发行所发行

2007 年 7 月第 1 版第 1 次印刷

184mm×260mm·26.75 印张

定价:49.00 元

凡购本书,如有倒页、脱页、缺页,由本社发行部调换
本社购书热线:(010)68326294

出版者的话

文艺复兴以降，源远流长的科学精神和逐步形成的学术规范，使西方国家在自然科学的各个领域中取得了垄断性的优势；也正是这样的传统，使美国在信息技术发展的六十多年间名家辈出、独领风骚。在商业化的进程中，美国的产业界与教育界越来越紧密地结合，计算机学科中的许多泰山北斗同时身处科研和教学的最前线，由此而产生的经典科学著作，不仅擘划了研究的范畴，还揭橥了学术的源变，既遵循学术规范，又自有学者个性，其价值并不会因年月的流逝而减退。

近年，在全球信息化大潮的推动下，我国的计算机产业发展迅猛，对专业人才的需求日益迫切。这对计算机教育界和出版界都既是机遇，也是挑战；而专业教材的建设在教育战略上显得举足轻重。在我国信息技术发展时间较短、从业人员较少的现状下，美国等发达国家在其计算机科学发展的几十年间积淀的经典教材仍有许多值得借鉴之处。因此，引进一批国外优秀计算机教材将对我国计算机教育事业的发展起积极的推动作用，也是与世界接轨、建设真正的世界一流大学的必由之路。

机械工业出版社华章图文信息有限公司较早意识到“出版要为教育服务”。自1998年开始，华章公司就将工作重点放在了遴选、移译国外优秀教材上。经过几年的不懈努力，我们与Prentice Hall, Addison-Wesley, McGraw-Hill, Morgan Kaufmann等世界著名出版公司建立了良好的合作关系，从它们现有的数百种教材中甄选出Tanenbaum, Stroustrup, Kernighan, Jim Gray等大师名家的一批经典作品，以“计算机科学丛书”为总称出版，供读者学习、研究及度藏。大理石纹理的封面，也正体现了这套丛书的品位和格调。

“计算机科学丛书”的出版工作得到了国内外学者的鼎力襄助，国内的专家不仅提供了中肯的选题指导，还不辞劳苦地担任了翻译和审校的工作；而原书的作者也相当关注其作品在中国的传播，有的还专程为其书的中译本作序。迄今，“计算机科学丛书”已经出版了近百个品种，这些书籍在读者中树立了良好的口碑，并被许多高校采用为正式教材和参考书籍，为进一步推广与发展打下了坚实的基础。

随着学科建设的初步完善和教材改革的逐渐深化，教育界对国外计算机教材的需求和应用都步入一个新的阶段。为此，华章公司将加大引进教材的力度，除“计算机科学丛书”之外，对影印版的教材，则单独开辟出“经典原版书库”。为了保证这两套丛书的权威性，同时也为了更好地为学校和老师服务，华章公司聘请了中国科学院、北京大学、清华大学、国防科技大学、复旦大学、上海交通大学、南京大学、浙江大学、中国科技大学、哈尔滨工业大学、西安交通大学、中国人民大学、北京航空航天大学、北京邮电大学、中山大学、解放军理工大学、郑州大学、湖北工学院、中国国家信息安全测评认证中心等国内重点大学和科研机构在计算机的各个领域的著名学者组成“专家指导委员会”，为我们提供选题意见和出版监督。

这两套丛书是响应教育部提出的使用外版教材的号召，为国内高校的计算机及相关专业的教学度身订造的。其中许多教材均已为M. I. T., Stanford, U. C. Berkeley, C. M. U.等世界名牌大学所采用。不仅涵盖了程序设计、数据结构、操作系统、计算机体系结构、数据库、编译原理、软件工程、图形学、通信与网络、离散数学等国内大学计算机专业普遍开设的核心课程，而且各具特色——有的出自语言设计者之手、有的历经三十年而不衰、有的已被全世界的几百所高

IV

校采用。在这些圆熟通博的名师大作的指引之下，读者必将在计算机科学的宫殿中由登堂而入室。

权威的作者、经典的教材、一流的译者、严格的审校、精细的编辑，这些因素使我们的图书有了质量的保证，但我们的目标是尽善尽美，而反馈的意见正是我们达到这一终极目标的重要帮助。教材的出版只是我们的后续服务的起点。华章公司欢迎老师和读者对我们的工作提出建议或给予指正，我们的联系方法如下：

电子邮件：hzsj@hzbook.com

联系电话：(010)68995264

联系地址：北京市西城区百万庄南街1号

邮政编码：100037

译者序

本书是算法设计与分析领域的经典之作。

本书和 Knuth 所著的《计算机程序设计艺术》^①、Cormen、Leiserson 和 Rivest 等著的《算法导论》^②合称为算法设计领域的 3 大名著。和《算法导论》力求浅显易懂及颇具全面性相比,本书更倾向于数学和形式化描述,特别对图灵机模型和空间复杂度等有更深入的讨论。本书的另一个特点是简洁,直指问题核心,并给出严谨的分析及解决方案和改进措施。

不管从内容上还是形式上,目前市场上所有算法设计的书籍都深受本书的影响(截至本书中文版出版之前,据 Amazon 统计,已有 271 本书籍引用了本书),可以说,不论是计算机理论研究人员、计算机算法设计者还是研究生、本科生,都可以通过阅读本书受到启发。

本书的特点

- 介绍了若干计算模型,包括图灵机模型、随机存取计算机模型以及它们的变体在算法设计和分析中的作用。
- 介绍了和设计有效算法相关的数据结构,给出递归、分治法、动态规划技术以及相关的算法例子。
- 解释了排序算法、集合算法、图算法、模式匹配算法、快速傅里叶变换的技术细节和应用技巧。
- 讨论了矩阵乘法、多项式运算以及整数算法的内在机理。
- 探讨了一些问题的时间复杂度和空间复杂度,涉及 NP 完全问题以及一些可以证明是不存在有效算法的问题,并将计算难度的概念和向量空间的线性独立性关联在一起。

关于本书附录

本书使用一种称为 Pidgin ALGOL 的语言来描述算法, Pidgin 意思是“由两种或多种语言混合而成的一种简化的说话方式,语法和词汇较初等,用于不同语言者进行交流”。算法的开发可分为设计、分析和实现 3 个环节。该过程可能不断重复,直到需求满足为止。一般说来,算法的描述只要清晰明了就可以了。书中所用的 Pidgin ALGOL 语言,对于算法的描述是充分的。但 ALGOL 是一个比较古老的语言,如果读者想测试书中给出的例子,那么找到一个 ALGOL 语言的编译器可是一件不容易的事(译者推荐使用 PLT Scheme 附带的 ALGOL 语言解释器,读者可参阅本书译者的另一本译作《程序设计方法》,人民邮电出版社 2003 年出版)。因此,在翻译本书的时候,经和机械工业出版社讨论协商,本书的代码保持原貌,另将书中的算法用 C/C++ 实现并作为译本的附录。本书附录中的代码由译者的研究生王慧芳实现,并经杨欢、彭冲、许赵云等同学测试,以保证其准确性。

本书的翻译工作由黄林鹏负责并与其博士生共同完成。其中第 9 章和第 10 章由王德俊翻译,第 7 章和第 8 章由张仕翻译,参加翻译的还有王欣、徐小辉、伍建焜等,全书由黄林鹏统稿和

① 该书第 1 卷第 1 册和第 4 卷第 2、3 和 4 册的双语版已由机械工业出版社出版。——编辑注

② 该书的第 2 版已由机械工业出版社翻译出版。——编辑注

VI

审校。

本书的英文版，曾学习过两次。一次是1980年，当时我在浙江大学读计算机本科，另一次是1986年，那时我已经在上海交通大学攻读硕士学位了。这本教材给我的印象很深，但没有想到在20多年后我有机会翻译这本书，感谢机械工业出版社华章分社给我重温巨著并用中文表述的机会。在此，我要向所有为本书的翻译提供帮助的人表示感谢。由于我们水平有限，翻译中难免出现不妥和错误之处，敬请读者谅解，并将意见寄至 lp Huang@sjtu.edu.cn，我们将不胜感激。

黄林鹏
于上海交通大学
2007年5月

前 言

算法研究是计算机科学的核心。近年来,算法领域取得了很多重要的进展。这些进展包括快速算法的开发,如发明了傅里叶变换快速算法,以及不存在有效算法的本质问题的惊人发现。这些结果点燃了计算机学者对算法研究的兴趣。算法设计与分析已成为一个受到广泛注意的领域。本书旨在将这个领域的一些基本成果汇集在一起,使算法设计的基本原则和根本原理的教学变得容易。

本书内容

分析一个算法的性能需要一种计算机模型。本书首先介绍一些形式模型,使用它们不仅可以进行算法分析,还能确切反映实际计算机的一些显著特性。这些模型包括随机存取计算机 RAM、随机存取存储程序计算机 RASP 以及一些变体。为了证明第 10、11 章中的一些算法的复杂度下界是指数的,本书还将讨论图灵机模型。由于程序设计趋向于远离机器语言,因此本书将引入一种称为简化 ALGOL 语言的高级程序语言作为描述算法的主要工具。简化 ALGOL 语言程序的复杂度和机器模型是相关的。

第 2 章介绍在设计有效的算法中使用的基本数据结构和程序设计技术。涉及列表、下推式存储结构、队列、树和图等。同时给出递归、分治法、动态规划的详细解释以及一些应用实例。

第 3~9 章给出一些不同的应用第 2 章所介绍的基本技术的领域。这些章节致力于设计已知最有效的渐近算法。由于要考虑渐近性,一些算法可能仅对那些输入规模比目前实际所遇到的问题更大时才有效。特别是第 6 章中的一些矩阵乘法算法、第 7 章中的 Schönhage - Strassen 整数乘法算法以及第 8 章中的一些多项式和整数算法。

另一方面,第 3 章中的大多数排序算法、第 4 章中的搜索算法、第 5 章的图算法、第 7 章的快速傅里叶算法以及第 9 章中的串匹配算法都是广泛应用的算法。这些算法对许多实际应用问题的输入规模来说都是有效的。

第 10~12 章讨论计算复杂度的下界。不管对于程序设计者还是希望了解计算本质的人都会对了解一个问题固有的计算难度感兴趣。第 10 章讨论一类重要的问题,即 NP 完全问题。该类中所有问题对于计算难度来说是等价的,即,若该类中有一个问题存在有效的算法(多项式时间界),则该类中的所有问题都有有效的算法。由于该问题类包含了许多在现实中非常重要且被广泛研究的问题,如整数规划问题和旅行商问题,有理由怀疑此类问题不存在有效的算法。因此,如果一个程序设计者了解到他欲寻找有效算法的问题属于此类,他就应该尝试寻找启发式的方法。尽管有丰富的经验证据,NP 完全问题是否存在有效的算法还是一个有待解决的问题。

第 11 章定义了一些确实可以证明不存在有效算法的问题。第 10 章和第 11 章的内容依赖于本书 1.6 和 1.7 节引入的图灵机的概念。

最后一章将计算难度的概念和向量空间的线性独立性关联在一起。该章给出了证明较第 10 章和第 11 章简单很多的一些问题的复杂度下界的方法。

如何使用本书

本书旨在作为涉及算法设计与分析课程的入门教材。强调算法的思想以及如何方便地理解

算法而不是算法实现的技巧和细节。因此，常使用内容直观的解释而不是冗长的证明。本书是自包含的，读者无需具有数学和程序设计语言的专业背景。然而，还是希望读者具备一定的处理数学概念的能力，熟悉某种高级程序设计语言如 FORTRAN 或 ALGOL 等。要完全理解第 6、7、8 和 12 章等内容，读者还需要具备线性代数的知识。

本书曾用于算法设计的本科生和研究生课程。一个学期的课程一般涉及第 1~5 章以及第 9 和 10 章的大部分内容，其他章节则可浅尝辄止。在算法的介绍性课程中，可以第 1~5 章为重点，而 1.6、1.7、4.13、5.11 节和定理 4.5 一般可不涉及。本书也可用于强调算法理论的高级课程，第 6~12 章的内容可作为此类课程的基础资料。

每章最后都提供有大量的习题，教师可选用作为学生的作业。习题按难度进行分类。没有星号的可用于介绍性的课程，有一个星号(*)的可用于高级课程，有两个星号的可用于研究生高级课程。每章后面的参考文献和注释提供了正文以及习题包含的相应算法和结果的公共资源。

致谢

本书的材料取自作者在康奈尔大学、普林斯顿大学和史迪温理工学院开设的算法课程的讲义。作者要感谢认真读过本书部分手稿的许多人所提出的批评和建议。作者特别要感谢 Kellogg Booth、Stan Brown、Steve Chen、Allen Cypher、Arch Davis、Mike Fischer、Hania Gajewska、Mike Garey、Udai Gupta、Mike Harrison、Matt Hecht、Harry Hunt、Dave Johnson、Marc Kaplan、Don Johnson、Steve Johnson、Brian Kernighan、Don Knuth、Richard Ladner、Anita LaSalle、Doug McIlroy、Albert Meyer、Christos Papadimitriou、Bill Plauger、John Savage、Howard Siegel、Ken Steiglitz、Larry Stockmeyer、Tom Szymanski 和 Theodore Yen 等。

特别感谢 Gemma Carnevale、Pauline Cameron、Hannah Kresse、Edith Purser 和 Ruth Suzuki 等对手稿所进行的细心和赏心悦目的排版。

感谢贝尔实验室、康奈尔大学、普林斯顿大学和加州大学伯克利分校为作者准备本书手稿所提供的设施。

A. V. A.

J. E. H.

J. D. U.

1974 年 6 月

目 录

出版者的话

译者序

前言

第 1 章 计算模型	1
1.1 算法和复杂度	1
1.2 随机存取计算机	3
1.3 RAM 程序的计算复杂度	7
1.4 存储程序模型	8
1.5 RAM 的抽象	11
1.6 一种基本的计算模型: 图灵机	15
1.7 图灵机模型和 RAM 模型的关系	19
1.8 简化 ALGOL——一种高级语言	20
第 2 章 有效算法的设计	26
2.1 数据结构: 表、队列和堆栈	26
2.2 集合的表示	28
2.3 图	29
2.4 树	30
2.5 递归	33
2.6 分治法	35
2.7 平衡	38
2.8 动态规划	39
2.9 后记	41
第 3 章 排序和顺序统计	46
3.1 排序问题	46
3.2 基数排序	47
3.3 比较排序	52
3.4 堆排序—— $O(n \log n)$ 的比较排序 算法	52
3.5 快速排序——期望时间为 $O(n \log n)$ 的排序算法	55
3.6 顺序统计学	58
3.7 顺序统计的期望时间	60
第 4 章 集合操作问题的数据结构	65
4.1 集合的基本操作	65
4.2 散列法	67

4.3 二分搜索	68
4.4 二叉查找树	69
4.5 最优二叉查找树	71
4.6 简单的不相交集合并算法	74
4.7 UNION-FIND 问题的树结构	77
4.8 UNION-FIND 算法的应用和扩展	83
4.9 平衡树方案	87
4.10 字典和优先队列	88
4.11 可合并堆	91
4.12 可连接队列	93
4.13 划分	94
4.14 本章小结	98
第 5 章 图算法	103
5.1 最小代价生成树	103
5.2 深度优先搜索	105
5.3 双连通性	107
5.4 有向图的深度优先搜索	112
5.5 强连通性	113
5.6 路径查找问题	117
5.7 传递闭包算法	119
5.8 最短路径算法	120
5.9 路径问题与矩阵乘法	121
5.10 单源问题	124
5.11 有向无环图的支配集: 概念 整合	126
第 6 章 矩阵乘法及相关操作	135
6.1 基础知识	135
6.2 Strassen 矩阵乘法算法	137
6.3 矩阵求逆	139
6.4 矩阵的 LUP 分解	140
6.5 LUP 分解的应用	145
6.6 布尔矩阵的乘法	146
第 7 章 快速傅里叶变换及其应用	153
7.1 离散傅里叶变换及其逆变换	153
7.2 快速傅里叶变换算法	156
7.3 使用位操作的 FFT	161
7.4 多项式乘积	164

7.5	Schönhage-Strassen 整数相乘 算法	165	10.2	\mathcal{P} 类和 \mathcal{NP} 类	231
第 8 章	整数与多项式计算	170	10.3	语言和问题	233
8.1	整数和多项式的相似性	170	10.4	可满足性问题的 NP 完全性	234
8.2	整数的乘法和除法	171	10.5	其他 NP 完全问题	239
8.3	多项式的乘法和除法	175	10.6	多项式空间界问题	245
8.4	模算术	177	第 11 章	一些可证难的问题	252
8.5	多项式模算术和多项式计值	179	11.1	复杂度层次	252
8.6	中国余数	180	11.2	确定型图灵机的空间层次	252
8.7	中国余数和多项式的插值	183	11.3	一个需要指数时间和空间的 问题	254
8.8	最大公因子和欧几里得算法	184	11.4	一个非基本的问题	260
8.9	多项式 GCD 的渐近快速算法	186	第 12 章	算术运算的下界	265
8.10	整数的 GCD	190	12.1	域	265
8.11	再论中国余数	191	12.2	再论直线状代码	266
8.12	稀疏多项式	192	12.3	问题的矩阵表述	267
第 9 章	模式匹配算法	196	12.4	面向行的矩阵乘法的下界	267
9.1	有穷自动机和正则表达式	196	12.5	面向列的矩阵乘法的下界	269
9.2	正则表达式的模式识别	201	12.6	面向行和列的矩阵乘法下界	272
9.3	子串识别	203	12.7	预处理	273
9.4	双向确定型下推自动机	207	附录	算法的 C/C++ 代码	280
9.5	位置树和子串标识符	215	参考文献	407
第 10 章	NP 完全问题	226			
10.1	非确定型图灵机问题	226			

第1章 计算模型

给定一个问题,如何寻求一个有效的求解算法?如果得到了一个算法又该如何与解决同一问题的其他算法进行比较?如何判断一个算法的好坏?程序设计者和理论计算机学者常对此类问题的本质感兴趣。本书将从不同的方面探讨这些问题。

本章将讨论几个典型的计算机模型,包括随机存取计算机 RAM、随机存取存储程序计算机 RASP 和图灵机 TM 等。比较这些模型的能力和算法复杂度的关系,从其导出若干更成熟的计算模型,如直线状程序模型、位计算模型、位向量计算模型和决策树模型等。本章的最后一节介绍一种描述算法的“简化 ALGOL 语言”。

1.1 算法和复杂度

评价算法有多个标准。最通常的是考虑其求解越来越大的问题实例所需的时间或空间开销的增长率。一般将一个称为问题规模的整数和算法相关联,该整数通常是问题的输入数据的规模。例如,对于矩阵乘法,其规模可能是相乘的矩阵的最大维数,而图问题的规模可能是边的数目。

算法所需要的时间是问题规模的函数,称为算法的时间复杂度。当问题规模增加时,复杂度的极限行为称为算法的渐近时间复杂度。类似的定义也适用于空间复杂度和渐近空间复杂度。

算法的渐近复杂度确定了可用该算法解决的问题的规模。对于某个常数 c , 如果一个算法能在 cn^2 时间内处理输入规模为 n 的问题,则称该算法的时间复杂度为 $O(n^2)$, 读为“阶为 n^2 ”。更精确地,如果存在常数 c , 除了有限的若干个值(可能为空)外,对所有非负的 n 都有 $g(n) \leq cf(n)$, 则称函数 $g(n)$ 是 $O(f(n))$ 的。

也许有人认为,现代数字计算机的出现带来的计算机速度的大幅度提高会削弱对有效算法的需求。然而,结果恰恰相反。计算速度的提高使我们可以处理更大的问题,但算法复杂度决定了由于计算速度的提高带来的可处理的问题规模的增量的大小。

假定有时间复杂度如下的 5 个算法 $A_1 \sim A_5$:

算法	时间复杂度
A_1	n
A_2	$n \log n^{\ominus}$
A_3	n^2
A_4	n^3
A_5	2^n

这里的时间复杂度是处理一个输入规模为 n 的问题的时间单元数。假定一个时间单元等于 1 毫秒,算法 A_1 能在一秒时间内处理输入规模为 1000 的问题,而算法 A_5 可在一秒内处理输入规模至多为 9 的问题。图 1-1 分别给出了 1 秒、1 分和 1 个小时内这 5 个算法可以处理的问题规模。

[⊖] 除非特别说明,本书所有对数都是以 2 为底。

算法	时间复杂度	最大问题规模		
		1 秒	1 分	1 小时
A_1	n	1000	6×10^4	3.6×10^6
A_2	$n \log n$	140	4893	2.0×10^5
A_3	n^2	31	244	1897
A_4	n^3	10	39	153
A_5	2^n	9	15	21

图 1-1 由增长率所确定的可解决的问题规模的限度

算法	时间复杂度	加速前最大问题规模	加速后最大问题规模
A_1	n	s_1	$10s_1$
A_2	$n \log n$	s_2	约 $10s_2$
A_3	n^2	s_3	对于大的 s_2 $3.16s_3$
A_4	n^3	s_4	$2.15s_4$
A_5	2^n	s_5	$s_5 + 3.3$

图 1-2 十倍加速的效果

假定下一代计算机的速度是目前的 10 倍。图 1-2 是计算速度增加后在相同的时间内可以解决的问题规模的增量。注意对于算法 A_5 ，10 倍的加速得到的结果是可以解决的问题规模比原来的规模仅大 3 左右，而对于算法 A_3 则是原来的 3 倍。

撇开速度的增加，现在考虑使用效率更高的算法的效果。回到图 1-1，使用 1 分钟为比较基础，用算法 A_3 代替算法 A_4 ，此时可以解决规模为原来 6 倍的问题；用算法 A_2 代替算法 A_4 ，则可以解决规模为原来 125 倍的问题。这些结果比用速度快 10 倍的计算机却得到 2 倍的改进更使人印象深刻。如果以 1 个小时作为比较基础，差异将更加明显。由此可以得到结论，一个算法的渐近时间复杂度是衡量算法好坏的一个重要标准，随着未来计算速度的提高，它必将变得更加重要。

除了注重阶量级外，也要认识到有时一个增长率较快的算法的比例常数可能比一个增长率较低的算法小。在这种情况下，对于小的问题，甚至是使我们感兴趣的同一规模下的所有问题，增长率较快的算法可能会比增长率低的算法优越。例如，假定算法 A_1 、 A_2 、 A_3 、 A_4 和 A_5 的时间复杂度分别是 $1000n$ ， $100n \log n$ ， $10n^2$ ， n^3 和 2^n ，则对于问题规模为 $2 \leq n \leq 9$ ， A_5 将是最好的算法；对于 $10 \leq n \leq 58$ ， A_3 将是最好的算法；对于 $59 \leq n \leq 1024$ ， A_2 将是最好的算法；而对于大于 1024 的规模， A_1 才是最好的算法。

在进一步讨论算法和复杂度之前，必须定义执行算法的计算设备模型，并指出计算中的基本步的含义。遗憾的是，并没有一个对所有情况都适用的模型。主要的困难在于如何对计算机字的大小进行定义。例如，如果假定一个计算机字可以存储任意大小的整数，则通过编码，整个问题就可以存放于一个计算机字中的一个整数。另一方面，如果假定一个计算机字的长度是有限的，则必须考虑如何存储任意大小的整数，以及通常在处理规模适当的问题时可以避免的其他困难。对于每个问题，都必须选择一个合适的模型来反映给定的算法在真正的计算机上运行时所耗费的时间。

接下来的几节讨论计算设备的一些基本模型，比较重要的模型是随机存取计算机、随机存取程序计算机和图灵机。在计算能力上它们是等价的，但在速度上不一样。

研究形式计算模型的最主要动机是揭示不同问题的固有计算难度。本章将证明一些问题的计算时间的下界。为了说明对于给定的问题不存在一个算法能在某一时间内解决它，需要一个精确的、常常是高度程式化的算法定义。图灵机(参见 1.6 节)就是定义算法的一个好的工具。

在描述算法和交流算法时，还需要一个比随机存取计算机、随机存取存储程序计算机或图灵机更加自然并易于理解的算法描述记法。由此，本书引进了称为 Pidgin(Pidgin 的含义是由两种或多种语言混合而成的一种简化的语言。——译者注)ALGOL 的高级语言，即简化 ALGOL 语言。本书将使用该语言进行算法描述。然而，为了解使用简化 ALGOL 语言描述的算法的复杂度，还必须将简化 ALGOL 语言和更加形式化的计算模型相关联。本章最后一节讨论该问题。

1.2 随机存取计算机

随机存取计算机(RAM)模拟一台带有一个累加器的计算机，该计算机的程序指令不允许对程序自身进行修改。

一台 RAM 计算机包括只读输入带、只写输出带、一个程序和一个存储器(见图 1-3)。输入带是一个方格序列，每个方格可存放一个整数(可以是负整数)。每从输入带读取一个整数，带首就往右移动一个方格。只写操作的输出带的初始状态全部为空。每执行一次写指令，就在当前输出带首的方格中写入一个整数，并将带首往右移动一个方格。不允许对输出带上写入的符号进行修改。

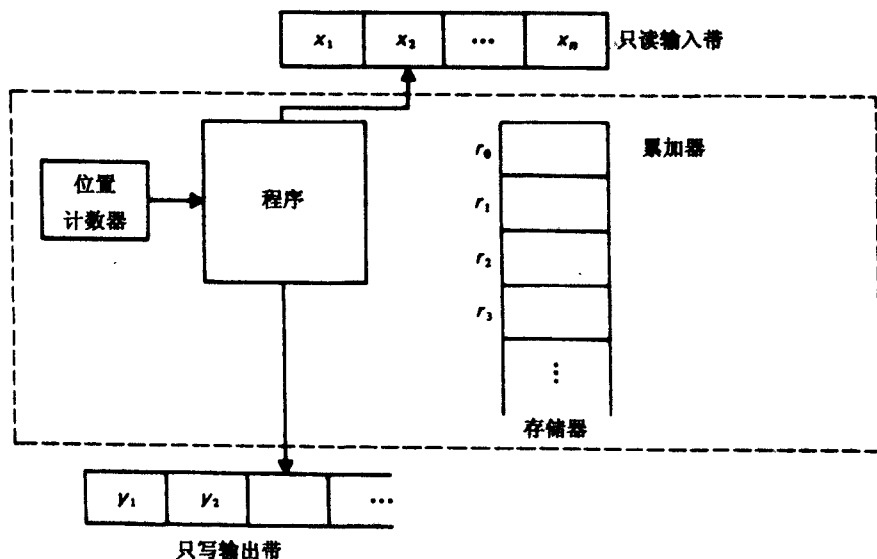


图 1-3 一台 RAM 计算机

存储器包含了寄存器 $r_0, r_1, \dots, r_i, \dots$ ，每个都可存储任意大小的整数，可用寄存器的数目没有限制。该抽象形式适用于下述情景：

1. 问题的规模足够小，可以存放在计算机的主存储器内；
2. 计算中用到的整数足够小，可以用一个计算机字表示。

在 RAM 模型中，程序并不存放在存储器中，因此程序不能修改自身。程序只是(可选)带标号的指令序列，同实际计算机上使用的指令类似，但其确切特性并不重要，通常可假定有算术运

算指令、输入输出指令、间接寻址(例如数组的下标)和转移指令等。所有的运算都在称为累加器的寄存器 r_0 中进行,累加器和其他的寄存器一样可以存储任意的整数。一个 RAM 指令集范例如图 1-4 所示。每条指令包含两个部分,操作码和地址。

原则上说,可以将图 1-4 中的指令集扩充为任何实际计算机的指令集,如添加逻辑或字符操作指令,而不改变解决问题的复杂度的阶量级。读者可以按照需要设想合适的指令集。

操作数可以是如下之一:

1. $=i$ 表示整数 i 本身;
2. 非负整数 i 表示寄存器 i 的内容;
3. $*i$ 表示间接寻址。即操作数是寄存器 j 的内容, j 是寄存器 i 中存放的整数。如果 $j < 0$, 则停机。

对于使用汇编语言编程的程序员来说,这些指令是再熟悉不过的了。定义程序 P 的含义,借助两个量,一个是从非负整数到整数的映射 c ,另一个是“位置计数器”,其确定下一条要执行的指令。函数 c 是内存映射, $c(i)$ 是存储在寄存器 i 中的整数(即寄存器 i 的内容)。

初始时,对所有 $i \geq 0$, $c(i) = 0$,位置计数器指向 P 中第一条指令,输出带全部为空白。在执行 P 的第 k 条指令之后,位置计数器的值自动设为 $k+1$ (即下一条指令),除非第 k 条指令是 JUMP、HALT、JGTZ 或 JZERO。

为了说明指令的含义,可定义操作数 a 的值 $v(a)$ 如下

$$\begin{aligned} v(=i) &= i \\ v(i) &= c(i) \\ v(*i) &= c(c(i)) \end{aligned}$$

图 1-5 中的表定义了图 1-4 中每条指令的含义。没有定义的指令,如 $\text{STORE } =i$, 可以认为和 HALT 一样。类似地,除以 0 也将使机器停机。

操作码	地址
1. LOAD	操作数
2. STORE	操作数
3. ADD	操作数
4. SUB	操作数
5. MULT	操作数
6. DIV	操作数
7. READ	操作数
8. WRITE	操作数
9. JUMP	标号
10. JGTZ	标号
11. JZERO	标号
12. HALT	

图 1-4 RAM 指令表

指令	含义
1. LOAD a	$c(0) \leftarrow v(a)$
2. STORE i	$c(i) \leftarrow c(0)$
STORE $*i$	$c(c(i)) \leftarrow c(0)$
3. ADD a	$c(0) \leftarrow c(0) + v(a)$
4. SUB a	$c(0) \leftarrow c(0) - v(a)$
5. MULT a	$c(0) \leftarrow c(0) \times v(a)$
6. DIV a	$c(0) \leftarrow \lfloor c(0)/v(a) \rfloor$ ①
7. READ i	$c(i) \leftarrow$ 当前输入符号
READ $*i$	$c(c(i)) \leftarrow$ 当前输入符号。这两种情况下,输入带首都右移一个方格
8. WRITE a	在输出带首的当前输出带方格打印 $v(a)$ 。然后带首右移一个方格
9. JUMP b	将位置计数器设置为标号为 b 的指令
10. JGTZ b	如果 $c(0) > 0$, 则将位置计数器设置为标号为 b 的指令; 否则, 将位置计数器的值设置为下一条指令
11. JZERO b	如果 $c(0) = 0$, 则将位置计数器设置为标号为 b 的指令; 否则, 将位置计数器设置为下一条指令
12. HALT	执行停止

① 在本书中, $\lceil x \rceil$ (x 的上取整) 表示大于或等于 x 的最小整数, $\lfloor x \rfloor$ (x 的下取整) 表示小于或等于 x 的最大整数。

图 1-5 RAM 指令的含义。操作数 a 是 $=i$ 、 i 或 $*i$

在执行头 8 条指令的每一条时, 位置计数器增 1。因此, 除非遇到 JUMP 或 HALT 指令, 或者当累加器的内容大于 0 时遇到 JGTZ 指令, 或者当累加器的内容等于 0 时遇到 JZERO 指令, 否则程序中的指令将按照顺序依次执行。

一般来说, RAM 程序定义了一个从输入带到输出带的映射。因为程序并不是对所有的输入带都会停机, 因此映射是一个部分映射(即, 对一些输入映射可能没有定义)。映射可以使用不同的方式解释, 其中两种重要的解释是把映射看作一个函数或一种语言。

假定程序 P 总是从输入带读取 n 个整数并最多在输出带上写进一个整数。如果 x_1, x_2, \dots, x_n 是输入带前 n 个方格中的整数, 程序 P 的执行结果是将 y 写入输出带的第一个方格并停机, 则称程序 P 计算了函数 $f(x_1, x_2, \dots, x_n) = y$ 。容易看出, RAM 模型和其他合理的计算机模型一样, 可以计算部分递归函数。即, 给定任意的部分递归函数 f , 可以定义一个计算 f 的 RAM 程序, 给定一个 RAM 程序, 可以定义一个等价的部分递归函数。(参阅 Davis[1958] 或 Rogers[1967] 对递归函数的讨论。)

另外一种解释是将 RAM 程序作为一个语言的接收器。字母表是有限个符号的集合, 语言是该字母表上字符串的一个集合。字母表的符号可以对于某个 k , 表示为整数 $1, 2, \dots, k$ 。RAM 可按以下方式接受一种语言。将输入串 $s = a_1 a_2 \dots a_n$ 放置到输入带, 其中符号 a_1 位于第一个方格、符号 a_2 位于第 2 个方格, 以此类推。另外, 将作为终止符号的 0 放置于第 $n+1$ 个方格, 作为输入符号串的开始。

如果 P 读入所有 s 中的符号以及结束符, 则输入串 s 被 RAM 程序 P 接受, 将 1 写入输出带的第 1 个方格并停机。 P 接受的语言是所有被 P 接受的输入符号串的集合。对不被 P 接受的输入字符串, 程序 P 将输出非 1 的符号并停机, 也可能不停机。容易证明, 一种语言能被 RAM 程序接受, 当且仅当它是递归可枚举的语言。一种语言能被对所有输入都可停机的 RAM 接受, 当且仅当它是递归语言(参见 Hopcroft 和 Ullman[1969] 对递归语言和递归可枚举语言的讨论)。

考虑下面两个 RAM 程序示例。前者定义一个函数, 后者接受一种语言。

例 1.1 考虑如下给定的函数 $f(n)$

$$f(n) = \begin{cases} n^n & \text{对所有整数 } n \geq 1 \\ 0 & \text{其他} \end{cases}$$

图 1-6 所示是一个简化 ALGOL 语言的程序^①, 通过将 n 自乘 $n-1$ 次来计算 $f(n)$ 。相应的 RAM 程序见图 1-7。其中变量 $r1$ 、 $r2$ 和 $r3$ 分别存放在寄存器 1、2 和 3 之中。由于没有进行程序优化, 图 1-6 和图 1-7 的对应比较明显。□

例 1.2 考虑一个 RAM 程序, 它接受输入字母表 $\{1, 2\}$ 上的语言, 该语言包含 1 和 2 的数目相等的任何字符串。程序将每个输入符号读入寄存器 1, 而寄存器 2 则保留了目前所看到的 1 和 2 的数目的差异数 d 。当读入终止符 0 时, 程序检查差异数是否为 0, 如是, 则打印 1 并停机。假定 0、1 和 2 是仅有的可能输入符号。

图 1-8 中的程序包含了基本的算法细节。等价的 RAM 程序如图 1-9 所示, 其中寄存器 1 存放 x , 而寄存器 2 存放 d 。□

```

begin
  read r1;
  if r1 ≤ 0 then write 0
  else
    begin
      r2 ← r1;
      r3 ← r1 - 1;
      while r3 > 0 do
        begin
          r2 ← r2 * r1;
          r3 ← r3 - 1
        end;
      write r2
    end
  end
end

```

图 1-6 计算 n^n 的简化 ALGOL 程序

① 参见 1.8 节对于简化 ALGOL 语言的描述。

RAM 程序		相应的简化 ALGOL 语句
	READ 1	read $r1$
	LOAD 1	
	JGTZ pos	if $r1 \leq 0$ then write 0
	WRITE =0	
	JUMP endif	
pos:	LOAD 1	$r2 \leftarrow r1$
	STORE 2	
	LOAD 1	$r3 \leftarrow r1 - 1$
	SUB =1	
	STORE 3	
while:	LOAD 3	while $r3 > 0$ do
	JGTZ continue	
	JUMP endwhile	
continue:	LOAD 2	$r2 \leftarrow r2 * r1$
	MULT 1	
	STORE 2	$r3 \leftarrow r3 - 1$
	LOAD 3	
	SUB =1	
	STORE 3	
	JUMP while	
endwhile:	WRITE 2	write $r2$
endif:	HALT	

图 1-7 计算 n^* 的 RAM 程序

```

begin
  d ← 0;
  read x;
  while x ≠ 0 do
    begin
      if x ≠ 1 then d ← d - 1 else d ← d + 1;
      read x
    end;
    if d = 0 then write 1
  end

```

图 1-8 识别具有相同数目的 1 和 2 的字符串的程序

RAM 程序		对应的混合 ALGOL 语句
	LOAD =0	$d \leftarrow 0$
	STORE 2	
	READ 1	read x
while:	LOAD 1	while $x \neq 0$ do
	JZERO endwhile	
	LOAD 1	if $x \neq 1$
	SUB =1	
	JZERO one	then $d \leftarrow d - 1$
	LOAD 2	
	SUB =1	
	STORE 2	
	JUMP endif	
one:	LOAD 2	else $d \leftarrow d + 1$
	ADD =1	
	STORE 2	
endif:	READ 1	read x
	JUMP while	
endwhile:	LOAD 2	if $d = 0$ then write 1
	JZERO output	
output:	HALT	
	WRITE =1	
	HALT	

图 1-9 与图 1-8 中算法相对应的 RAM 程序

1.3 RAM 程序的计算复杂度

衡量算法的两个重要标准是时间和空间复杂度，它们都是输入规模的函数。如果对于给定的输入规模，复杂度考虑的是对所有给定规模的输入，取其中最大者，则该复杂度称为最坏情况复杂度；而若复杂度考虑的是对所有给定规模的输入，取“平均”值，则该复杂度称为期望复杂度。算法的期望复杂度通常比最坏情况复杂度更难以确定，这是因为前者通常必须假定输入的分布情况，而现实的输入分布在数学上并不容易处理。因此，本书将注重最坏情况复杂度，它既容易处理也在实际中有着广泛的应用。然而，必须记住的是算法在最坏情况复杂度最好的算法并不必然是期望时间复杂度最好的算法。

一个 RAM 程序的最坏情况时间复杂度（或简称时间复杂度）是函数 $f(n)$ ，表示对规模为 n 的所有输入，所有被执行的指令的“时间”之和的最大值。期望时间复杂度则是对规模为 n 的所有输入，所有被执行的指令的“时间”之和的平均值。若将“执行指令花费的时间”换为“寄存器占用的空间”，就得到了空间复杂度。

精确定义时间和空间复杂度还必须规定执行每条 RAM 指令所花费的时间和每个寄存器使用的空间。对 RAM 程序，这里将考虑两种代价标准。在统一代价标准中，每条 RAM 指令的执行花费一个单位的时间，而每个寄存器占用一个单位的空间。本书如果没有特别指出，RAM 程序的复杂度就按统一代价标准进行度量。

有时另一衡量标准更现实，它考虑到实际存储器字的有限大小，称为对数代价标准。令 $l(i)$ 是如下定义在整数上的对数函数：

$$l(i) = \begin{cases} \lfloor \log |i| \rfloor + 1 & i \neq 0 \\ 1 & i = 0 \end{cases}$$

如图 1-10 所示的表给出了对于操作数 a 的 3 种可能形式，对数代价 $t(a)$ 。图 1-11 归纳了每条指令所需的时间开销。

操作数 a	代价 $t(a)$
$=i$	$l(i)$
i	$l(i) + l(c(i))$
$*i$	$l(i) + l(c(i)) + l(c(c(i)))$

图 1-10 操作数的对数代价

指令	代价
1. LOAD a	$t(a)$
2. STORE i	$l(c(0)) + l(i)$
STORE $*i$	$l(c(0)) + l(i) + l(c(i))$
3. ADD a	$l(c(0)) + t(a)$
4. SUB a	$l(c(0)) + t(a)$
5. MULT a	$l(c(0)) + t(a)$
6. DIV a	$l(c(0)) + t(a)$
7. READ i	$l(\text{input}) + l(i)$
READ $*i$	$l(\text{input}) + l(i) + l(c(i))$
8. WRITE a	$t(a)$
9. JUMP b	1
10. JGTZ b	$l(c(0))$
11. JZERO b	$l(c(0))$
12. HALT	1

图 1-11 RAM 指令的对数代价，这里 $t(a)$ 表示操作数 a 的开销， b 是标号

这里考虑的事实是，在一个寄存器中表示整数 n 需要 $\lfloor \log n \rfloor + 1$ 位，而前面说一个寄存器可以存放任意大小的整数。

对数代价模型基于如下粗略假设，执行一条指令的代价和指令的操作数的长度成比例。例如，考虑指令 ADD $*i$ 的代价。首先必须考虑对操作数的地址进行解码的开销。检查整数 i 需要时间 $l(i)$ 。然后，读出 $c(i)$ ，即寄存器 i 的值，再定位寄存器 $c(i)$ ，需要时间 $l(c(i))$ 。最后读取寄存器 $c(i)$ 的内容，需要时间 $l(c(c(i)))$ 。由于指令 ADD $*i$ 是将 $c(c(i))$ 的值和累加器中的

整数 $c(0)$ 相加, 可以看出 $l(c(0)) + l(i) + l(c(i)) + l(c(c(i)))$ 是指令 $\text{ADD } *i$ 的合理的时间开销。

定义一个 RAM 程序的对数空间复杂度是 $l(x_i)$ 值的和, 这里 x_i 是包括累加器在内的所有寄存器 i 在计算过程中所存放过的最大整数。

对于给定的程序使用统一代价标准或对数代价标准可能会得到非常不同的时间复杂度分析结果。如果假定程序中的每个数都可存放在一个计算机字内, 则统一代价函数是合理的模型。否则, 使用对数代价模型进行实际复杂度分析更合适。

现在考虑例 1.1 中计算 n^n 的 RAM 程序的时间和空间复杂度。该程序的时间复杂度由包含 MULT 指令的循环所决定。执行第 i 次 MULT 指令时, 累加器中包含 n^i 的值, 而寄存器 2 则包含值 n 。一共执行了 $n-1$ 次 MULT 指令。在统一代价标准下, 每条 MULT 花费一个单元的时间, 因此执行所有的 MULT 指令花费 $O(n)$ 时间。在对数代价模型下, 执行第 i 条 MULT 指令花费的时间是 $l(n^i) + l(n) \approx (i+1) \log n$, 因此 MULT 总的执行时间是

$$\sum_{i=1}^{n-1} (i+1) \log n$$

即 $O(n^2 \log n)$ 。

空间复杂度由存储于寄存器 0~3 中的整数所确定。在统一代价标准下空间复杂度是简单的 $O(1)$ 。在对数代价标准下, 空间复杂度是 $O(n \log n)$ 。因此对于例 1.1, 程序的复杂度如下:

	统一代价	对数代价
时间复杂度	$O(n)$	$O(n^2 \log n)$
空间复杂度	$O(1)$	$O(n \log n)$

对于该程序, 仅当一个计算机字能存放 n^n 这样大的整数时, 使用统一代价标准才是现实的。如果 n^n 无法存放在一个计算机字内, 由于假定整数 i 和 j 相乘需时间 $O(l(i) + l(j))$, 而这是是否可能尚不清楚, 因此有时甚至使用对数代价标准来衡量算法的时间复杂度也是不实际的。

对于例 1.2 中的 RAM 程序, 假定 n 是输入串的长度, 其时间和空间复杂度如下:

	统一代价	对数代价
时间复杂度	$O(n)$	$O(n \log n)$
空间复杂度	$O(1)$	$O(n \log n)$

对于这个程序, 如果 n 不能存放在一个计算机字内, 则对数代价标准更现实。

1.4 存储程序模型

由于 RAM 程序并不存放在 RAM 的内存中, 因此程序不能对自身进行修改。现在考虑一种称为随机存取存储程序计算机 (RASP) 的模型, 它和 RAM 模型类似, 唯一的不同点是程序存放在内存之中并且可以修改。

RASP 的指令集和 RAM 基本一样, 不同之处是没有间接寻址指令, 原因是不需要。RASP 可以在程序执行过程中修改指令来仿真间接寻址。

RASP 的总体结构也和 RAM 相似, 但 RASP 的程序存放在内存的寄存器之中。每条 RASP 指令占用两个连续的内存寄存器。第 1 个寄存器存放表示操作码的代码, 第 2 个寄存器存放地址。如果地址的形式是 $=i$, 则第 1 个寄存器存放的代码将表示第 2 个寄存器的内容是 i 。使用整数对指令进行编码。图 1-12 给出一种可能的编码。例如, 指令 $\text{LOAD } =32$ 将在一个寄存器中存放

值 2 而在紧接的寄存器中存放值 32。

指令	编码	指令	编码
LOAD i	1	DIV i	10
LOAD $=i$	2	DIV $=i$	11
STORE i	3	READ i	12
ADD i	4	WRITE i	13
ADD $=i$	5	WRITE $=i$	14
SUB i	6	JUMP i	15
SUB $=i$	7	JGTZ i	16
MULT i	8	JZERO i	17
MULT $=i$	9	HALT	18

图 1-12 RASP 指令编码

相应于 RAM, RASP 的状态可表示如下:

1. 内存映射 c , 对所有 $i \geq 0$, $c(i)$ 为寄存器 i 的内容;
2. 位置计数器, 其指向两个连续寄存器的第 1 个寄存器, 从该处取当前要执行的指令。

开始时, 位置计数器设置为指向某一指定的寄存器。初始时, 因为有装入的程序, 内存寄存器的内容并不全是 0。但可假定, 除了有限个寄存器外, 其他寄存器和累加器的值都是 0。除非指令是 JUMPi 、 JGTZi (且累加器值是正数) 或 JZEROi (且累加器值为 0), 此时位置计数器的值设为 i , 否则每条指令执行后, 位置计数器增加 2。每条指令的效果和相应的 RAM 指令一样。

RASP 程序的时间复杂度的定义非常类似于 RAM。既可以使用统一代价标准也可以使用对数代价标准。但对于后者, 必须考虑操作数的计算以及存取指令本身的代价。存取指令的代价是 $l(\text{LC})$, 这里 LC 是位置计数器的值。例如, 执行存放在寄存器 j 和 $j+1$ 中的指令 $\text{ADD} = i$ 的代价是 $l(j) + l(c(0)) + l(i)$ 。^①

执行存放在寄存器 j 和 $j+1$ 中的指令 $\text{ADD} i$ 的代价是 $l(j) + l(c(0)) + l(i) + l(c(i))$ 。

RAM 程序的计算复杂度和 RASP 程序的计算复杂度的不同之处在哪里? 答案并不令人意外。一个模型中时间复杂度为 $T(n)$ 的从输入到输出的映射, 在另一个模型中在时间 $kT(n)$ 内完成 (k 为常数), 无论统一代价还是对数代价均如此。类似地, 在这两种代价模型下, 这两个模型所耗用的空间也仅差常数倍。

下面两个定理形式化地阐述了上述关系。定理的证明是通过说明一个 RAM 算法可以由一个 RASP 算法进行仿真, 反之亦然。

定理 1.1 不管指令代价是统一的或者是对数的, 对于每一个时间复杂度为 $T(n)$ 的 RAM 程序, 都存在一个常数 k , 使得存在一个等价的时间复杂度为 $kT(n)$ 的 RASP 程序。

证明: 下面说明如何用 RASP 程序来仿真一个 RAM 程序 P 。RASP 的寄存器 1 用来暂时存储 RAM 累加器的内容。现在由 P 构造一个 RASP 程序 P_s , P 占用 RASP 接下来的 $r-1$ 个寄存器。常数 r 由 RAM 程序 P 决定。RAM 寄存器 i 的内容存放在 RASP 寄存器 $r+i$ 中, 其中 $i \geq 1$, 由此所有对 RASP 程序的内存引用都比对相应的 RAM 程序的内存引用多 r 。

每一条没有包含间接引用的 RAM 程序 P 的指令可直接编码为相同的 RASP 指令 (内存引用可能适当增加)。每一条包含间接引用的 RAM 程序 P 的指令都可映射为通过更改指令以仿真间接寻址的 6 条 RASP 指令序列。

① 也可考虑读取寄存器 $j+1$ 的代价, 但和 $l(j)$ 没有大的不同。另外本章不关心常数因子, 而是考虑函数的增长率。因此 $l(j) + l(j+1)$ “近似” 于 $l(j)$, 即相差最多 3 倍。

下面的例子足以说明如何进行间接寻址仿真。为了仿真 RAM 指令 $\text{SUB } *i$, 其中 i 是一个正整数, 现编译 RASP 指令序列如下:

1. 将累加器中的内容临时存放在寄存器 1 中;
2. 将寄存器 $r+i$ 的内容装入累加器中(RASP 模型的寄存器 $r+i$ 对应于 RAM 模型的寄存器 i);
3. 将累加器内容加 r ;
4. 将第 3 步得到的结果作为 SUB 指令的地址域存储;
5. 从临时寄存器 1 恢复累加器的值;
6. 最后, 使用第 4 步创建的 SUB 指令执行减法运算。

寄存器	内容	含义
100	3	STORE 1
101	1	
102	1	
103	$r+i$	LOAD $r+i$
104	5	
105	r	
106	3	ADD $=r$
107	111	
108	1	
109	1	STORE 111
110	6	
111	-	

图 1-13 通过 RASP 模型仿真指令 $\text{SUB } *i$

例如, 使用图 1-12 给出的 RASP 指令编码并假定 RASP 指令系列开始于寄存器 100, 图 1-13 所示为 $\text{SUB } *i$ 的仿真指令序列。偏移量 r 在得到 RASP 程序 P_s 后就可以确定。

可以观察到每条 RAM 指令最多需要 6 条 RASP 指令, 因此在统一代价标准下 RASP 指令的复杂度最多为 $6T(n)$ 。(注意该度量与确定输入“规模”的方法无关。)

在对数代价模型下, P 中的每条 RAM 指令 I 可以被 P_s 下 1 条或 6 条指令序列 S 仿真。可以看到存在一个依赖于 P 的常数 k , 使得 S 中的指令代价不大于指令 I 的代价的 k 倍。

例如, RAM 指令 $\text{SUB } *i$ 的代价为

$$M = l(c(0)) + l(i) + l(c(i)) + l(c(c(i)))$$

仿真该 RAM 指令的 RASP 指令序列 S 如图 1-14 所示, 其中 $c(0)$ 、 $c(i)$ 和 $c(c(i))$ 是 RAM 寄存器的内容。因为 P_s 占用 RASP 的寄存器 $2 \sim r$, 由此, $j \leq r-11$ 。由于 $l(x+y) \leq l(x) + l(y)$, 因此, S 的代价肯定小于

$$2l(1) + 4M + 11l(r) < (6 + 11l(r))M$$

RASP寄存器	指令	代价
j	STORE 1	$l(j) + l(1) + l(c(0))$
$j+2$	LOAD $r+i$	$l(j+2) + l(r+i) + l(c(i))$
$j+4$	ADD $=r$	$l(j+4) + l(c(i)) + l(r)$
$j+6$	STORE $j+11$	$l(j+6) + l(j+11) + l(c(i) + r)$
$j+8$	LOAD 1	$l(j+8) + l(1) + l(c(0))$
$j+10$	SUB -	$l(j+10) + l(c(i) + r) + l(c(0)) + l(c(c(i)))$

图 1-14 RASP 指令的代价

由此可得结论, 存在一个常数 $k = 6 + 11l(r)$, 使得若 P 的时间复杂度为 $T(n)$, 则 P_s 的时间复杂度为 $kT(n)$ 。□

定理 1.2 不管指令代价模型是统一的或者是指数的, 对于每一个时间复杂度为 $T(n)$ 的 RASP 程序, 都存在常数 k , 使得存在一个等价的时间复杂度为 $kT(n)$ 的 RAM 程序。

证明: 要构造的仿真 RASP 的 RAM 程序将使用间接寻址以解码并仿真存放在 RAM 内存中的 RASP 程序。其中一些 RAM 的寄存器有特殊的使用目的:

寄存器 1——用作间接寻址,

寄存器 2——RASP 的位置计数器,

寄存器 3——存储用作 RASP 的累加器。

RASP 的寄存器 i 存放在 RAM 的寄存器 $i+3$ 之中, $i \geq 1$ 。

RAM 将长度有限的 RASP 程序存放在开始于寄存器 4 的内存中。寄存器 2, 即位置计数器的值为 4, 而寄存器 1 和 3 的值为 0。RAM 程序包含一个仿真循环, 其中首先读取一条 RASP 指令 (LOAD *2 的 RAM 指令), 然后译码, 转移到 18 种不同的指令集之一, 每个指令集处理一种类型的 RASP 指令。对于无效的 RAM 操作码, 和 RASP 类似, 将停机。

译码和转移操作是非常直接的, 例 1.2 可以作为一个模型 (虽然此例中解码的符号是来自输入, 而现在是从内存读取)。图 1-15 是仿真 RASP 指令 6, 即 SUB i 的 RAM 指令序列, 当 $c(c(2)) = 6$, 即位置计数器指向一个含有 6 (SUB 的编码) 的寄存器时该程序被激活。

LOAD	2	将位置计数器的值增 1, 使其指向包含 SUB i 指令的操作数 i 的
ADD	=1	寄存器
STORE	2	
LOAD	*2	将 i 存入累加器, 加 3, 再将结果存放在寄存器 1 中
ADD	=3	
STORE	1	
LOAD	3	从寄存器 3 中得到 RASP 累加器的内容, 减去寄存器 $i+3$ 的值, 将
SUB	*1	结果放回寄存器 3 中
STORE	3	
LOAD	2	
ADD	=1	再次将位置计数器的值增 1, 使其指向下一条 RASP 指令
STORE	2	
JUMP	a	返回到循环仿真指令序列的开始处 (假定是 “ a ”)

图 1-15 仿真 SUB i 的 RAM 代码

这里省略掉 RAM 程序构造的进一步细节。请读者自行证明: 不管是统一代价模型还是对数代价模型, RAM 程序的时间复杂度最多是 RASP 程序的常数倍。□

从定理 1.1、定理 1.2 可以得出以下结论, 对于时间复杂度 (空间复杂度也是, 证明留作练习), RAM 模型和 RASP 模型是常数因子等价的, 也就是说相同的算法在两个模型下复杂度的阶是一样的。涉及这两个模型时, 本书通常将使用较简单的 RAM 模型。

1.5 RAM 的抽象

对于许多情况, RAM 和 RASP 比问题所需要的计算模型复杂。下面将定义若干对 RAM 模型的一些属性进行抽象而忽略其他特性的模型。考虑准则是所忽略的指令的代价最多占应用该模型的任何有效算法的代价的常数部分。

1. 直线状程序

第一个考虑的模型是直线状程序模型。对许多问题来说, 有理由仅考虑如下 RAM 程序, 其中转移指令仅用于表示一些指令序列的执行会重复 n 次, 这里 n 是问题的输入规模。在这种情况下

下, 可以对重复 n 次的指令序列进行适当次数的复制而“展开”程序, 其结果就是长度可能增加 n 倍的直线状程序(无循环)。

例 1.3 考虑两个 $n \times n$ 的整数矩阵相乘。可以认为在 RAM 程序中一个循环被执行的次数独立于实际的矩阵元素的取值。因此, 假定程序中出现循环测试条件语句仅涉及 n , 即问题的规模, 这是一个有用的简化。例如, 矩阵乘法算法显然涉及恰好执行 n 次的循环, 即其包含和 n 进行比较的转移指令。□

将一个程序展开为直线状程序可以省却分支指令。展开的理由是在许多问题中, 用于控制循环的分支指令的代价总和只是整个 RAM 程序代价的常量因子。同样, 通常也可假设输入语句也是程序总代价的一个常量因子, 因此可以假设大小为 n 的有限输入集在程序开始时就已经存放在内存中。另外, 如果用于间接寻址的寄存器包含的值仅依赖于 n , 而与输入的变量值无关, 那么间接寻址的影响当 n 固定时也可以确定。由此, 可以假定直线状程序不包含间接寻址指令。

另外, 由于每一直线状程序仅引用有限个内存寄存器, 由此可方便地对程序使用的寄存器命名, 即可使用符号地址(符号或字符串)而不是整数来引用寄存器。

消去 READ、JUMP、JGTZ、JZERO 之后, 所剩下的指令还包括 LOAD、STORE、WRITE、HALT 以及 RAM 指令系统中的算术运算等。由于程序的结束必然是程序的停止, 因此 HALT 指令也是不需要的。另外, 还能通过指定某一符号地址作为输出变量而免去 WRITE 指令, 此时在程序停止时, 这些地址持有的值就是程序的输出。

最后, 通过使用下面的指令序列可将 LOAD 和 STORE 指令和算术指令结合在一起, 如将

```
LOAD    a
ADD     b
STORE   c
```

用 $c \leftarrow a + b$ 替换。由此, 直线状程序的整个指令集为:

```
x ← y + z
x ← y - z
z ← y * z
z ← y / z
x ← i
```

这里 x 、 y 和 z 是符号地址(或变量), 而 i 是一个常量。容易看出, 累加器中由 LOAD、STORE 和算术运算组成的指令序列都可替换为上述 5 条指令组成的指令序列。

和直线状程序相关的是两组指定的变量, *inputs* 和 *outputs*。直线状程序所计算的是输入变量表示的值到输出变量(按指定顺序)表示的值的函数。

例 1.4 考虑多项式求值

$$p(x) = a_n x^n + a_{n-1} x^{n-1} + \cdots + a_1 x + a_0$$

输入变量是系数 a_0, a_1, \cdots, a_n 和未确定的 x 。输出变量是 p 。计算 $p(x)$ 的霍纳规则为

1. $a_1 x + a_0$ 对于 $n=1$
2. $(a_2 x + a_1)x + a_0$ 对于 $n=2$
3. $((a_3 x + a_2)x + a_1)x + a_0$ 对于 $n=3$

图 1-16 是与这些表达式对应的直线状程序。对于任何 n , 霍纳规则现在应该是清楚的。对于任何一个 n , 都有步数为 $2n$ 的计算阶为 n 的多项式的直线状程序。第 12 章将说明对于给定的输入系数计算一个任意的 n 阶多项式, n 次加法和 n 次乘法运算是必需的。因此对于直线状程序模型, 使用霍纳规则的程序是最优的。□

$n = 1$	$n = 2$	$n = 3$
$t \leftarrow a_1 * x$	$t \leftarrow a_2 * x$	$t \leftarrow a_3 * x$
$p \leftarrow t + a_0$	$t \leftarrow t + a_1$	$t \leftarrow t + a_2$
	$t \leftarrow t * x$	$t \leftarrow t * x$
	$p \leftarrow t + a_0$	$t \leftarrow t + a_1$
		$t \leftarrow t * x$
		$p \leftarrow t + a_0$

图 1-16 相应于霍纳规则的直线状程序

在直线状程序计算模型下,一段程序序列的时间复杂度是第 n 个程序的计算步数,即 n 的函数。例如,霍纳规则产生的是时间复杂度为 $2n$ 的一个序列。注意,测量时间复杂度和测量算术运算的数目的结果是一样的。程序序列的空间复杂度是涉及到的变量的数目,它也是 n 的函数。例 1.4 中的程序空间复杂度为 $n+4$ 。

定义 当涉及直线状程序模型时,对于一个问题,若存在一个时间或空间复杂度最多为 $cf(n)$ 的程序, c 为常数,则称该问题的时间或空间复杂度为 $O_A(f(n))$ 。(记法 $O_A(f(n))$ 表示“使用直线状程序模型的阶为 $f(n)$ ”。下标 A 表示“算术”,是直线状程序代码的主要特征。)因此,计算多项式的时间复杂度为 $O_A(n)$,空间复杂度也为 $O_A(n)$ 。

II. 按位计算

直线状程序模型显然基于统一代价函数。正如前面所提及的,该代价模型假设所有计算的量的大小“合理”时是合适的。对直线状程序模型进行简单的修改就可得到对数代价模型。该模型可称为按位计算模型,除了下面一些性质外和直线状程序模型是一样的:

1. 所有变量的值为 0 或 1,即变量是位变量;
2. 运算是逻辑操作而不是算术操作。^①

用 \wedge 表示 and, \vee 表示 or, \oplus 表示异或, \neg 表示否定。

在按位计算模型中,整数 i 和整数 j 的算术运算至少需要 $l(i) + l(j)$ 步,反映了操作数的对数代价。事实上,已知最好的乘法和除法运算的算法确实需要多于 $l(i) + l(j)$ 步来计算 i 和 j 的积和商。

使用 O_b 来表示按位计算模型的算法复杂度的阶。当涉及在其他模型中的基本运算,如算术操作时,按位计算模型是有用的。例如,在直线状程序模型下,两个 n 位整数相乘可在 $O_A(1)$ 步内完成,而在按位计算模型下已知的最好的结果是 $O_b(n \log n \log \log n)$ 步。

按位计算模型的另一个应用领域是逻辑电路。输入和运算为二进制的直线状程序和计算布尔函数的组合逻辑电路存在一一对应的关系。程序中计算的步数就是电路中的逻辑单元数。

例 1.5 图 1-17a 是将两个二进制 2 位数 $[a_1 a_0]$ 和 $[b_1 b_0]$ 相加的程序。输出变量是 c_2 、 c_1 和 c_0 ,使得 $[a_1 a_0] + [b_1 b_0] = [c_2 c_1 c_0]$ 。图 1-17a 的直线状程序计算为:

$$\begin{aligned} c_0 &= a_0 \oplus b_0 \\ c_1 &= ((a_0 \wedge b_0) \oplus a_1) \oplus b_1 \\ c_2 &= ((a_0 \wedge b_0) \wedge (a_1 \vee b_1)) \vee (a_1 \wedge b_1) \end{aligned}$$

图 1-17b 是对应的逻辑电路。作为练习,读者请自行证明可在 $O_b(n)$ 步内完成两个 n 位数相加的操作。 □

① 因此 RAM 的指令集必须包含这些操作。

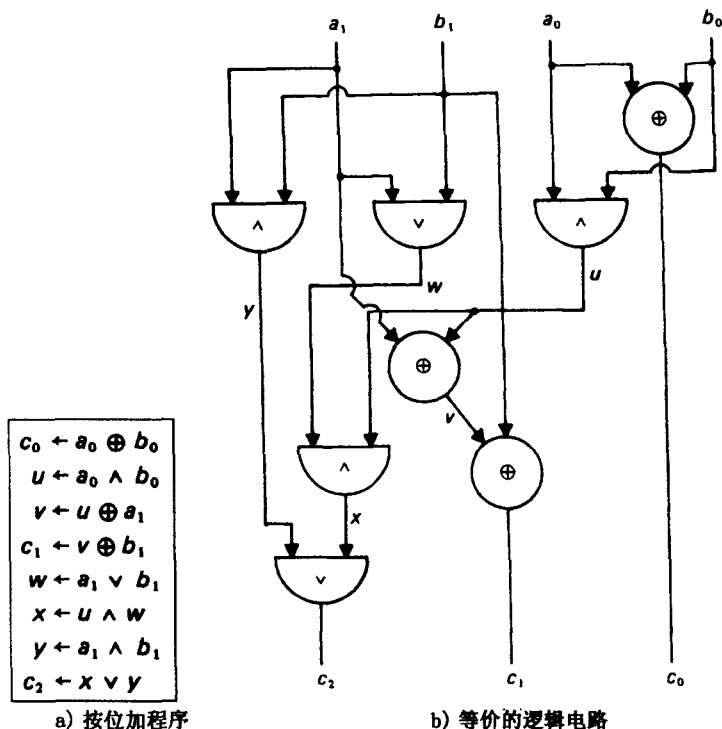


图 1-17

III. 位向量操作

与将变量的值限制为 0 和 1 相反的方法是允许变量的值可以是任意长度的位向量。实际上，固定长度的位向量显然对应于整数，因此，必须采用比 RAM 模型更本质的自由度，即方便时可假设寄存器的大小是无界的。

虽然有一些算法在使用位向量模型时，所用向量的长度会超出表示问题规模所用的位数，但大部分算法中使用的整数的阶和问题规模的阶是一样的。例如，处理一个有 100 个顶点的图的路径问题，可使用长度为 100 的位向量来表示一个给定的顶点 v 是否和其他的顶点间存在一条路径，即当且仅当存在从 v 到 v_i 的一条路径，则与顶点 v 相关的向量的第 i 位值为 1。对于同一问题，还可以使用规模大小是 100 的整数作为计数和索引，此时，需要使用 7 位的整数，而对于向量则需要 100 位。

然而，比较结果并不全是一边倒的。大多数计算机在一个指令周期内只能完成一个字长度的位向量逻辑运算。因此，相应于一步完成的整数操作，长度为 100 的位向量需要 3~4 步内完成。因为模型变得非现实时，问题的规模比 RAM 和直线状代码模型更小，因此，应用位向量模型的算法的时间和空间复杂度有时是有疑问的。本书使用 O_{bv} 表示使用位向量模型的阶。

IV. 决策树

前面已经讨论了 3 种仅考虑计算步骤而忽略分支转移指令的 RAM 抽象，但也存在一些问题。有时分支转移指令是度量复杂度的主要因素。例如，在排序问题中，输出除了次序外和输入是一样的。因此有理由考虑这样的模型，全部代码都是基于两个元素比较的两分支转移指令。

程序中表示分支的通常工具是使用称为决策树的二叉树。[⊖]每个内部顶点表示一个决策。树

⊖ 树的定义请参阅 2.4 节。

根所表示的测试首先进行, 然后按照测试的结果将“控制”转移给它的一个儿子, 通常, 依赖于结点测试结果。控制持续从一个顶点转移到它的一个子孙, 直到叶子。所期望的结果就在所抵达的叶子之上。

例 1.6 图 1-18 说明了对 3 个数 a 、 b 和 c 进行排序的决策树程序。结点处圆圈内的比较式表示所进行的测试, 如果结果是“yes”, 则控制转移到左儿子, 结果是“no”, 则控制转移到右儿子。□

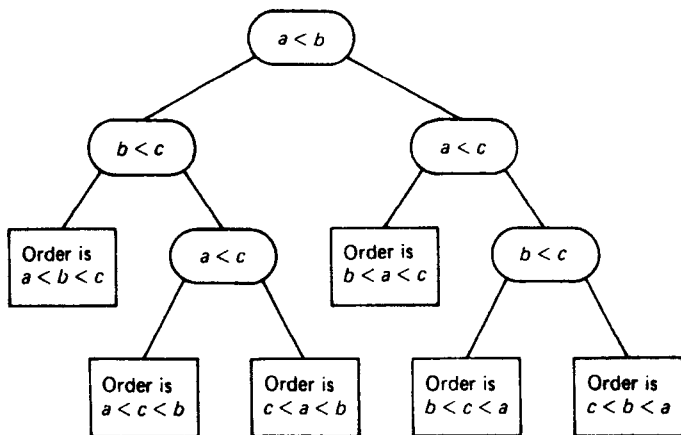


图 1-18 一棵决策树

决策树的时间复杂度是树的高度, 是问题规模的函数。通常, 我们希望测量的是最大的比较次数, 是从根到叶子的路径长度。本书用 O_c 表示在决策树(比较)模型下算法的阶。注意, 一棵树的结点的总数通常超过它的高度。例如, 对 n 个数进行排序的决策树至少有 $n!$ 片叶子, 尽管其高度大约 $n \log n$ 就足够了。

1.6 一种基本的计算模型: 图灵机

为了证明计算一个特殊函数需要不可避免的最小时间开销, 我们需要一种一般的、但比前面所讨论的更加基本的计算模型。该模型的指令集要尽可能少, 同时它不仅能计算 RAM 模型所能计算的一切函数, 而且要“差不多”一样快。“差不多”的定义可使用“多项式相关”概念。

定义 若存在多项式 $p_1(x)$ 和 $p_2(x)$, 使对于所有 n 值, 有 $f_1(n) \leq p_1(f_2(n))$ 和 $f_2(n) \leq p_2(f_1(n))$, 则说函数 $f_1(n)$ 和 $f_2(n)$ 是多项式相关的。

例 1.7 函数 $f_1(n) = 2n^2$ 和 $f_2(n) = n^5$ 是多项式相关的; 可令 $p_1(x) = 2x$, 因为 $2n^2 \leq 2n^5$ 。令 $p_2(x) = x^3$, 因为 $n^5 \leq (2n^2)^3$ 。然而, n^2 和 2^n 不是多项式相关的, 因为不存在多项式 $p(x)$, 使对于所有 n 有 $p(n^2) \geq 2^n$ 。□

目前, 使用如图灵机这样的通用计算模型来证明计算复杂性的下界是在“较宽的范围”进行的。例如, 第 11 章说明某些问题需要指数阶的时间和空间。 $(f(n))$ 是一个指数函数, 如果存在常数 $c_1 > 0$ 、 $k_1 > 0$ 、 $c_2 > 0$ 和 $k_2 > 1$, 对除有限个例外的所有 n , 有 $c_1 k_1^n \leq f(n) \leq c_2 k_2^n$ 。在一个指数的范围, 多项式相关函数本质上是一样的, 因为和一个指数函数是多项式相关的函数本身必定是一个指数函数。

因此有动机使用一个基本的模型, 在该模型上问题的时间复杂度和该问题在 RAM 模型上的时间复杂度是多项式相关的。具体来说, 使用的模型是多带图灵机, 以对数代价函数计, 需要

($[f(n)]^4$) 时间完成一台 RAM 机器在 $f(n)$ 时间内完成的工作, 但不会更多^①。因此 RAM 和图灵机模型的时间复杂度是多项式相关的。

定义 一台多带图灵机(TM)如图 1-19 所示。它包含 k 条右端无限的带子, 每条带子划分为格子, 每个格子可以包含一个带符号, 取自一个有限符号集。磁头可扫描每条带子的一个格子, 进行读或写。图灵机的操作是由一个称为有限控制器的基本程序确定的。有限控制器总是处于一个有限状态集中的某一状态, 该状态可以被认为是程序的一个位置。

下面是图灵机上的一个计算步骤。基于有限控制器的当前状态以及每个磁头所扫描的相应符号, 图灵机执行下述操作中的一些或全部:

1. 改变有限控制器的状态。
2. 打印新的带符号以覆盖若干或全部磁头下的格子的当前符号。
3. 独立移动若干或全部磁头, 往左(L)一格、往右(R)一格或不变(S)。

形式化的描述是, 一台 k 带图灵机是 7 元组

$$(Q, T, I, \delta, b, q_0, q_f)$$

其中:

1. Q 是状态的集合;
2. T 是带符号的集合;
3. I 是输入符号的集合, $I \subseteq T$;
4. b 属于 $T - I$, 是空白符;
5. q_0 是初始状态;
6. q_f 是最终(接受)状态;
7. δ 为移动函数, 是从 $Q \times T^k$ 到 $Q \times (T \times \{L, R, S\})^k$ 的映射。即, 对包含一个状态和 k 个带符号的 $(k+1)$ 元组, 给出一个新的状态和 k 个对, 每个对包含一个新的带符号和磁头的移动方向。假定 $\delta(q, a_1, a_2, \dots, a_k) = (q', (a_1', d_1), (a_2', d_2), \dots, (a_k', d_k))$, 而图灵机在状态 q , 并将第 i 条磁带的磁头扫描带符号是 a_i , $1 \leq i \leq k$ 。接下来, 图灵机进入状态 q' , 将符号 a_i 改为 a_i' , 并将第 i 条磁带的磁头按照方向 d_i 移动, $1 \leq i \leq k$ 。

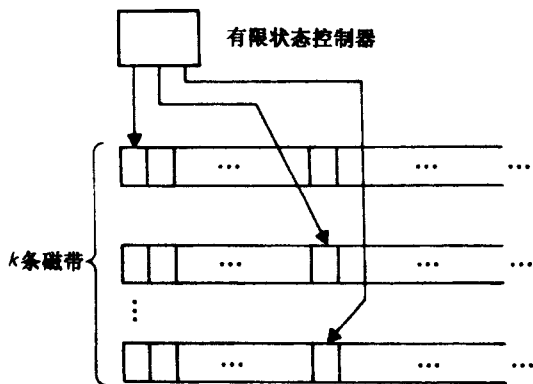


图 1-19 一台多带图灵机

① 事实上, 可以证明一个更紧密的界 $O([f(n) \log f(n) \log \log f(n)]^2)$, 但由于这里关心的是多项式因素, 使用 4 次方的界就可以了(见 7.5 节)。

可以使用一台图灵机来识别一种语言。图灵机的带符号包括输入语言的字母表,称为输入符号,和一个特殊符号,即记为 b 的空白符,或许还有其他的符号。开始时,第一条带上包含输入符号串,从最左边开始每个格子一个符号。包含输入符号串格子的右边的所有格子都是空白。其他的带子全为空白符。当且仅当图灵机从指定的初始状态,即每一磁头都处于磁带最左的位置出发,在进行一系列的移动之后进入可接受状态,称输入符号串是可接受的。一台图灵机接受的语言是所有可接受的输入符号串的集合。

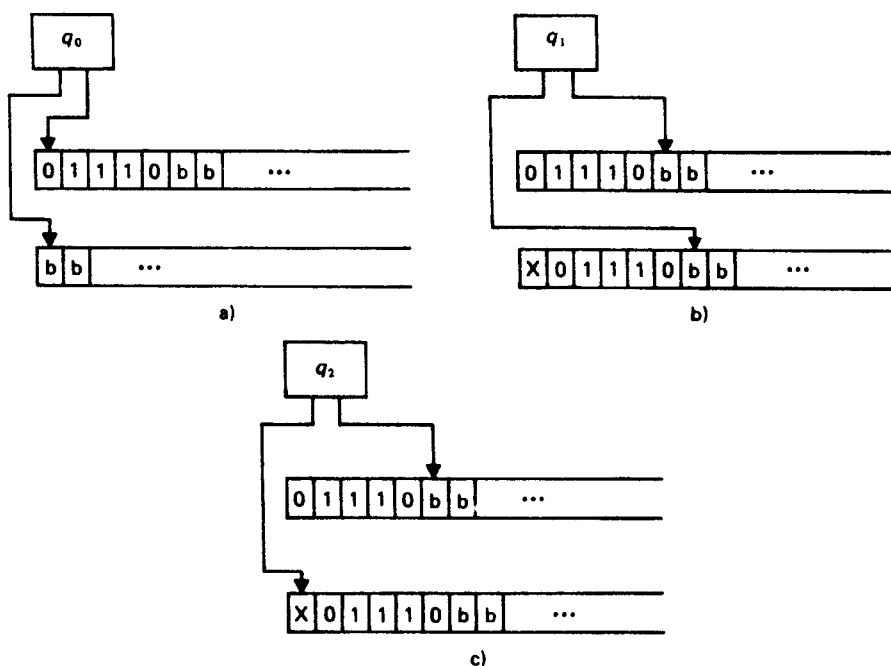


图 1-20 处理 01110 的图灵机

例 1.8 图 1-20 是识别字母表 $\{0, 1\}$ 上回文^①的 2 带图灵机。

1. 在第 2 条带的第 1 个格子内写入特殊的符号 X, 并从第 1 条带复制初始的输入串(见图 1-20a)到第 2 条带(见图 1-20b);
2. 将第 2 条带的磁头移到符号 X 所在的位置(见图 1-20c);
3. 重复如下动作, 磁带 2 上的磁头每次向右移动 1 格同时磁带 1 上的磁头向左移动 1 格, 并比较相应磁头读出的符号。如果所有的符号都是匹配的, 则可判定输入串是一回文, 此时图灵机进入状态 q_5 。否则, 图灵机将位于一个无法进行合法移动的位置, 此时将不接受字符串而停机。

图灵机的下一移动函数由图 1-21 中的表给出。

□

图灵机的动作可以使用“瞬像”进行形式化描述。一台 k 带图灵机 M 的瞬像(instantaneous description, ID)是一个 k 元组 $(\alpha_1, \alpha_2, \dots, \alpha_k)$, 其中 α_i 是形为 xqy 的串, 其中 xy 是 M 中第 i 条带上的串(略去尾上的空字), q 是 M 的当前状态。仅挨着第 i 个 q 的右边符号就是第 i 条带扫描的符号。

① 一字符串若从前往后读与从后往前读一样则称为回文, 如 0100010。

当前状态	符号		新符号, 磁头移动		新状态	注 解
	带 1	带 2	带 1	带 2		
q_0	0	b	0,S	X,R	q_1	如果输入非空, 在带 2 上打印 X 并将磁头右移, 进入状态 q_1 ; 否则进入状态 q_5
	1	b	1,S	X,R	q_1	
	b	b	b,S	b,S	q_5	
q_1	0	b	0,R	0,R	q_1	停留在状态 q_1 , 将带 1 的内容复制到带 2 直至遇到带 1 上的符号 b。然后进入状态 q_2
	1	b	1,R	1,R	q_1	
	b	b	b,S	b,L	q_2	
q_2	b	0	b,S	0,L	q_2	带 1 上的磁头不动, 带 2 上的磁头向左移动直至遇到符号 X。然后进入状态 q_3
	b	1	b,S	1,L	q_2	
	b	X	b,L	X,R	q_3	
q_3	0	0	0,S	0,R	q_4	控制器在状态 q_3 和 q_4 之间变换。在 q_3 , 比较两条带上的符号, 将带 2 的磁头右移, 并进入状态 q_4 。在 q_4 , 如果带 2 上的磁头遇到 b, 进入状态 q_5 并接受输入串; 否则左移带 1 的磁头并回到状态 q_3 。状态 q_3 和 q_4 之间的变换防止了带 1 上的磁头从最左边滑出
	1	1	1,S	1,R	q_4	
q_4	0	0	0,L	0,S	q_3	
	0	1	0,L	1,S	q_3	
	1	0	1,L	0,S	q_3	
	1	1	1,L	1,S	q_3	
	0	b	0,S	b,S	q_5	
	1	b	1,S	b,S	q_5	
q_5						接受

图 1-21 识别回文的图灵机的下一移动函数

如果经图灵机 M 一次移动后瞬像 D_1 转成瞬像 D_2 , 则记为 $D_1 \vdash D_2$ (读为“转成”, 或“进入”)。若 $D_1 \vdash D_2 \vdash \cdots \vdash D_n$, 对于某个 $n \geq 2$, 则记为 $D_1 \vdash^* D_n$ 。若 $D = D'$ 或 $D \vdash^* D'$, 则记为 $D \vdash^* D'$ 。

k 带图灵机 $M = (Q, T, I, \delta, b, q_0, q_f)$ 接受串 $a_1 a_2 \cdots a_n$, 其中 a_i 属于 I , $1 \leq i \leq n$, 如果 $(q_0 a_1 a_2 \cdots a_n, q_0, q_0, \cdots, q_0) \vdash^* (\alpha_1, \alpha_2, \cdots, \alpha_k)$, 其中某个 α_i 包含 q_f 。

例 1.9 图 1-22 是图 1-21 表示的图灵机对于输入 010 的瞬像序列。由于 q_5 是最终状态, 因此该图灵机接受 010。

除了作为语言接收器的自然解释之外, 图灵机还被认为是用于计算函数 f 的设备。函数的参数编码为输入带上的符号串 x , 使用特殊的符号如 # 来分开不同的参数。如果图灵机停机时, 作为输出的带上有整数 y (函数的值), 则说 $f(x) = y$ 。因此, 计算一个函数的过程和接受一个语言的过程有稍许不同。

图灵机 M 的时间复杂度 $T(n)$ 是 M 在处理任何长度为 n 的输入时计算的最大步数。如果对于某个长度为 n 的输入, 图灵机不停机, 则 $T(n)$ 对该值无定义。图灵机的空间复杂度 $S(n)$ 是在处理长度为 n 的输入时, 磁头离开带子左端的最大距离。如果有一条带子的磁头往右无限移动, 则 $S(n)$ 无定义。本书使用 O_{TM} 表示图灵机模型下的复杂度。

例 1.10 当输入确实是一串回文时, 可以验证图 1-21 中图灵机的时间复杂度 $T(n) = 4n +$

$(q_0 010, q_0)$	$(q_1 010, X q_1)$
	$(0 q_1 10, X 0 q_1)$
	$(01 q_1 0, X 01 q_1)$
	$(010 q_1, X 010 q_1)$
	$(010 q_2, X 010 q_2 0)$
	$(010 q_2, X 0 q_2 10)$
	$(010 q_2, X q_2 010)$
	$(010 q_2, q_2 X 010)$
	$(01 q_2 0, X q_2 010)$
	$(01 q_2 0, X 0 q_2 10)$
	$(0 q_2 10, X 0 q_2 10)$
	$(0 q_2 10, X 01 q_2 0)$
	$(q_2 010, X 01 q_2 0)$
	$(q_2 010, X 010 q_2)$
	$(q_2 010, X 010 q_2)$
	$(q_5 010, X 010 q_5)$

图 1-22 图灵机的 ID 序列

3, 空间复杂度 $S(n) = n + 2$ 。

□

1.7 图灵机模型和 RAM 模型的关系

图灵机(TM)模型的主要应用是确定求解给定问题所需的时间和空间的下界。大多数情况下, 所能确定的是多项式相关范围的下界。要推出更紧密的下界则需要具有更明确细节的模型。幸运的是, 在 RAM 或 RASP 上的计算和 TM 上的计算是多项式相关的。

本节考虑 RAM 和 TM 模型的相关性。显然, 让 RAM 的每个寄存器包括 TM 带子上每个格子的内容, RAM 能仿真一台 k 带 TM。特别是, 第 j 条带子的第 i 个格子的内容可存放在第 $ki + j + c$ 寄存器上, 其中 c 是一个常量, 表示 RAM“工作空间”。工作空间中有 k 个用于存放 TM 的 k 个磁头的位置信息的寄存器。TM 带子上的格子内容可由 RAM 通过引用存有磁头位置的寄存器使用间接寻址进行读取。

假定 TM 的时间复杂度为 $T(n) \geq n$ 。这样 RAM 可以读取输入, 然后将输入存放在模拟第 1 条磁带的寄存器上, 若使用统一代价函数就可在 $O(T(n))$ 内仿真 TM 计算, 若使用对数代价模型则可在 $O(T(n)\log T(n))$ 内仿真 TM 计算。不管如何, RAM 程序计算的时间开销的上界是和 TM 上计算的开销是多项式相关的, 因为任何 $O(T(n)\log T(n))$ 函数肯定是 $O(T^2(n))$ 的。

对于逆向仿真, 仅对使用对数代价函数的 RAM 模型是有效的。在统一代价模型下, 一个没有输入的 n 步 RAM 程序可以得到高达 2^n 的整数值, 需要 2^n 个 TM 格子来存放和读取。因此, 在统一代价下, RAM 和 TM 之间不存在多项式相关性是显然的(习题 1.19)。

虽然由于简单性, 在分析算法时倾向于使用统一代价函数, 但当想证明算法时间复杂度的下界时则不行。使用统一代价的 RAM 模型当涉及的数没有过大而超过问题规模比例时是比较合理的。但正如前面提到的, RAM 模型下数的大小像被“扫到地毯下的垃圾”, 很少能得到有用的下界。而对于对数代价, 则有下面的定理。

定理 1.3 令 L 是使用对数代价标准的 RAM 程序在时间复杂度 $T(n)$ 内接受的语言。如果 RAM 程序不使用乘法或除法, 则最多花费 $O(T^2(n))$ 时间, L 就可被一台多带图灵机接受。

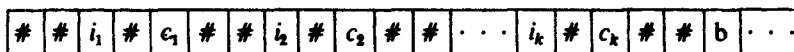


图 1-23 RAM 的 TM 表示

证明: 将所有内容不为 0 的寄存器表示为图 1-23 所示的带 1, 该带包括 (i_j, c_j) 的序列, 用抹去高位端的 0, 并用符号“#”隔开的二进制表示。对于每个 j , c_j 是 RAM 寄存器 i_j 的内容。RAM 累加器的内容以二进制的形式存放在第 2 条带上, 第 3 条带则用作工作存储。另外两条带则用于存放 RAM 输入和输出的内容。RAM 程序的每一步表示为 TM 的一个有限状态集。本书并不打算描述对任意的 RAM 指令的仿真, 仅考虑指令 ADD * 20 和 STORE 30 的仿真, 目的是理清技术思路。对于 ADD * 20, 可设计能完成如下工作的 TM:

1. 在带 1 上查找相应于 RAM 寄存器 20 的位置, 即序列##10100#。如果找到, 则将其后的整数, 即寄存器 20 的内容, 写到带 3。如果没有找到, 则停机。若寄存器 20 的内容为 0, 则无法进行间接寻址。
2. 在带 1 上查找是否有 RAM 寄存器其编号存储在带 3 上。如果找到, 则将其内容复制到带 3。如果没有找到, 则写入 0。
3. 将第 2 步放置在带 3 上的数和带 2 上的累加器中的内容相加。

为了仿真指令 STORE 30, 可设计能完成如下工作的 TM:

1. 在带 1 上查找相应于 RAM 寄存器 30 的位置, 即##11110#。

2. 如果找到, 将##11110#右边的所有符号, 除了直接紧随的整数(寄存器 30 的旧内容)以外, 复制到带 3。然后将累加器的内容(带 2)直接复制到##11110#的右边, 再将带 3 的内容写回。
3. 如果在带 1 上没找到寄存器 30, 则磁头移动至最左的空白处, 打印 11110#, 接着打印累加器的内容和分隔符“##”。

稍微想一下就可以确信, 可以使用一台 TM 来忠实仿真 RAM 模型, 但还必须说明一个对数代价为 k 的 RAM 计算最多需要图灵机上 $O(k^2)$ 步的计算。注意如下现象, 除非一个寄存器的当前值在先前某一时间存于其内, 否则该寄存器是不会有带 1 中出现的。将 c_j 存放于寄存器 i_j 的代价是 $l(c_j) + l(i_j)$, 在一个常数的误差下, 这和表示## i_j ## c_j ##的长度是一样的。由此, 可以得到带 1 的非空白符部分的长度为 $O(k)$ 。

因为主要的开销是对带的搜索, 因此任何非 STORE 指令的仿真开销和带 1 的长度是同一阶的, 即 $O(k)$ 。同样, 一条 STORE 指令的代价最多为对带 1 的搜索代价加上复制该带内容的代价, 两者都是 $O(k)$ 。因此一条 RAM 指令(除了乘法和除法指令)可以在 TM 中使用至多 $O(k)$ 步进行仿真。因为一条 RAM 指令在对数代价下至少花费一个单位的时间, 因此可以证明 TM 的总花费时间是 $O(k^2)$ 。□

如果一个 RAM 程序采用乘法和除法指令, 那么可以编写使用加法和减法指令的 TM 子例程来实现它们。读者可以自行证明在对数代价下这些子例程的代价不多于它们所仿真的指令的对数代价的平方。因此, 不难证明如下定理。

定理 1.4 在对数代价下, RAM 模型、RASP 模型和多带图灵机模型是多项式相关的模型。

证明: 使用定理 1.1、定理 1.2 和定理 1.3 以及对乘法和除法子例程的分析。□

对于空间复杂度也有同样的结论, 但结果没有那么有意思。

1.8 简化 ALGOL——一种高级语言

尽管计算复杂度可借助于 RAM、RASP 或图灵机上的操作步数进行度量, 但一般不使用也不必用如此基本的机器指令来描述算法。为了更清晰地描述算法, 本书将使用一种高级语言, 即简化 ALGOL 语言。

一个简化 ALGOL 语言程序可以直接转化为 RAM 或 RASP 程序, 更精确地说, 这实际上是简化 ALGOL 语言编译器的任务。然而, 本书并不关心将简化 ALGOL 语言的程序转化为 RAM 代码或 RASP 代码的细节。本书的目的是关注执行简化 ALGOL 语句所需的时间和空间开销。

和其他传统程序设计语言不同, 简化 ALGOL 语言将使用任何类型的数学表达式, 只要表达式含义清楚并且容易翻译为 RAM 或 RASP 代码即可。类似地, 该语言也没有固定的数据类型集。变量可以是整型、字符串类型和数组类型。额外的数据类型如集合、图、表和队列等, 也可以在需要的时候引入。本书会尽可能避免对数据类型进行形式说明。一个变量的数据类型和其作用域。^①可以从其名字或上下文识别出来。

简化 ALGOL 使用传统的数学和程序设计语言结构, 如表达式、条件、语句和过程等。下面是一些此类结构的非形式化描述。对它们的精确定义超出了本书的范围。要认识到, 虽然没有细节的精确说明, 还是可以写出程序, 因此, 尽可能忽略无关的细节, 本书(希望)如此进行。

① 变量的作用域是一个环境, 在其中该变量有意义。例如, 一个求和算符的指标的作用域仅仅在求和表达式中有定义, 在求和表达式外是没有意义的。

简化 ALGOL 程序包括如下类型的语句。

1. **variable** \leftarrow **expression**
2. **if condition then statement else statement**[⊖]
- 3a. **while condition do statement**
 - b. **repeat statement until condition**
4. **for variable** \leftarrow **initial-value step step-size until final-value do statement**[⊖]
5. **label: statement**
6. **goto label**
7. **begin**
 - statement;**
 - statement;**
 - .**
 - .**
 - statement;**
 - statement**
- end**
- 8a. **procedure name (list of parameters): statement**
 - b. **return expression**
 - c. **procedure-name (arguments)**
- 9a. **read variable**
 - b. **write expression**
10. **comment comment**
11. 其他形式的语句

下面是这些语句类型的一个简要说明。

1. 赋值语句

variable \leftarrow **expression**

计算 \leftarrow 右边表达式的值，并将结果赋予 \leftarrow 左边的变量。赋值语句的时间复杂度是计算表达式的时间加上将值赋予变量的时间。如果表达式的值不是如整数这样的基本类型，那么一些情况下可以使用指针进行简化。例如，对于 $n \times n$ 矩阵 A 和 B ，赋值语句 $A \leftarrow B$ 的开销通常是 $O(n^2)$ 。然而，如果 B 不再使用，则只要使用重命名技术就可以将时间控制在与 n 无关的有限步内。

2. if 语句

if condition then statement else statement

if 之后的条件可以是值为 **true** 或 **false** 的任意表达式。如果条件式为 **true**，则 **then** 后的语句将被计算。否则 **else** 后的语句（如果有的话）会被计算。if 语句的开销是计算和测试条件式的开销加上计算实际被执行的 **then** 之后的语句或 **else** 之后的语句的开销。

3. while 语句

while condition do statement

和 **repeat** 语句

repeat statement until condition

⊖ “**else statement**”是可选的，该选项导致所谓的“悬垂 **else**”歧义性，一般采用传统的方法，即假定 **else** 和最近的未匹配的 **then** 相匹配。

⊖ 当 **step-size** 为 1 时，“**step step-size**”是可选的。

两种语句的目的是创建循环。在 **while** 语句中, **while** 后的条件式先计算。如果条件式为 **true**, 则 **do** 后的语句被执行。该过程一直重复直到条件变为 **false** 为止。如果条件开始时为 **true**, 要使 **while** 语句的执行最后停止, 语句的执行最终会让条件成为 **false**。**while** 语句的代价是对条件多次求值的代价的总和, 加上语句多次被执行的代价的总和。

除了紧随 **repeat** 之后的语句在条件表达式之前被计算外, **repeat** 语句的计算与 **while** 语句类似。

4. **for** 语句

for variable ← initial-value step step-size until final-value do statement

initial-value、**step-size** 和 **final-value** 都是表达式。当 **step-size** 为正数时, **variable** (称为索引) 的值设为初始表达式的值。如果该值超过了 **final-value**, 执行停止。否则 **do** 之后的语句被执行, **variable** 的值增加 **step-size** 并再次和 **final-value** 进行比较。该过程一直进行直到 **variable** 的值超过了 **final-value**。**step-size** 是负数的情况与此类似, 当 **variable** 的值小于 **final-value** 时执行停止。借鉴前面对于 **while** 语句的分析, **for** 语句的执行代价是显而易见的。

上述分析完全忽略了表达式 **initial-value**、**step-size** 和 **final-value** 的求值是何时进行的细节。**do** 后面的语句的执行很可能改变表达式 **step-size** 的值, 在这种情况下每次对 **step-size** 求值还是计算一次以后一直使用是有区别的。类似地, 对 **step-size** 的求值可能会影响 **final-value** 的值, 同样 **step-size** 的改变也会影响终止测试的结果。为避免这些问题, 本书将避免编写这些含义不清的程序。

5. 任何一条语句都可以在其前面加上标号和冒号成为带标号的语句。标号的主要目的是创建 **goto** 语句的转移目标。标号没有附加的代价。

6. **goto** 语句

goto label

使带有该标号的语句成为下一执行的语句。与标号相关的语句不允许出现在一个语句块内部(见下面的第7点), 除非 **goto** 语句出现在同一语句块内。**goto** 语句的代价为1。由于 **goto** 语句会带来对程序的理解困难, 所以应该节俭使用。**goto** 语句的主要目的是跳出一个 **while** 语句。

7. 由分号隔开并包含在关键字 **begin** 和 **end** 之间的若干语句称为块(block)。由于块是一个语句, 它可以出现在任何语句出现的地方。通常, 一个程序是一个块。块的代价是出现在块内的语句代价的总和。

8. 过程。在简化 ALGOL 中, 过程可以先定义然后调用。过程通过形式如下的过程定义语句定义

procedure name(list of parameters): statement

其中 **list of parameters** 是称为形式参数的亚元变量的序列。例如, 下列语句定义了一个名为 **MIN** 的函数过程

procedure MIN(x, y):
if x > y then return y else return x

其中变量 x 和 y 是形式参数。

有两种过程形式。一是函数, 定义一个函数过程之后, 就可以在一个表达式内使用函数名字和相应的参数调用它。在这种情况下, 过程的最后一条被执行的语句必须是 **return** 语句, 即 8(b)。**return** 语句的执行使关键字 **return** 之后的表达式被计算, 并使过程的执行停止。函数的值就是该表达式的值。例如

$A \leftarrow \text{MIN}(2 + 3, 7)$

使 A 的值为5。表达式 $2 + 3$ 和 7 称为调用过程的实际参数。

另一种使用过程的方法是通过过程调用语句 8(c)。该语句形式是过程名后跟随一个实参表。

过程调用语句(通常)可以改变调用程序的数据。按这种方式调用过程无需在过程定义体内存在 **return** 语句。过程中最后一条语句的执行使过程调用语句的执行结束。例如, 下面的语句定义了一个名为 **INTERCHANGE** 的过程

```

procedure INTERCHANGE (x, y):
begin
    t  $\leftarrow$  x;
    x  $\leftarrow$  y;
    y  $\leftarrow$  t
end

```

为了调用该过程, 编写如下的过程调用语句

```
INTERCHANGE (A[i], A[j])
```

一个过程有两种和其他过程进行通信的方式。一种称为全局变量。假定全局变量是在全局 (universal) 环境下隐式说明的。而过程是在全局环境下的一个子环境中定义的。

另外一种通信方式是使用参数。ALGOL60 使用按值调用和按名调用。在按值调用下, 将过程的形式参数被处理为初始值为实参值的局部变量。在按名调用下, 形式参数被看成是程序中的一个位置占位符, 形式参数的每一个出现将被实参替换。为了简化起见, 不使用 ALGOL60 的机制, 而是采用按引用调用的机制。在该机制下, 参数作为指向实际参数的指针传递。如果实际参数是一个表达式(可能是常数), 则对应的形式参数可处理为初始值为表达式值的局部变量。在 RAM 和 RASP 中函数或过程调用的代价是执行过程定义体中的语句的代价总和。第 2 章将讨论调用其他过程(包括自身调用)的执行代价。

9. **read** 语句和 **write** 语句的执行代价是明显的。**read** 语句的代价为 1。**write** 语句的代价是 1 加上对关键字 **write** 后面的表达式求值的代价。

10. **comment** 语句允许在程序中插入有助于程序理解的注解, 其代价为 0

11. 除了通常程序设计语言的语句外, 还包括一些“其他形式”的语句, 与等价的程序语言语句序列相比, 这些语句可提高算法的易读性。这些语句通常在实现细节与算法无关、实现细节明显或需要高层次的描述时使用。通常使用的其他形式的语句的例子包括:

- a) 令 a 是集合 S 中的最小元素;
- b) 标记元素 a 为“旧的”[⊖]
- c) 不失一般性 (**wlg**) 假定...否则...**in** 语句

例如,

wlg assume $a \leq b$ **otherwise** interchange c and d **in** statement

意思是, 如果 $a \leq b$, 则直接执行后续的语句; 如果 $a > b$, 则在后续的语句执行过程中 c 和 d 的角色互换。

使用通常的程序设计语言语句或 RAM 代码实现这些语句直接而乏味。这些语句的执行代价依赖于语句所在的上下文。在本书的简化 ALGOL 程序中可以找到采用这种类型的语句的实例。

因为通常不对变量进行说明, 所以这里将介绍一下关于变量作用域的使用惯例。在一个给定的程序或过程中, 不同的变量不使用同一个名字。因此, 可以将变量的作用域看作是变量出现在其中的整个过程或程序。[⊖]重要的例外是可能存在一个共用的数据块, 有几个过程共同对其操

⊖ 可假定有一个数组 **STATUS**, 若 a 是“旧的”, 则 **STATUS**[a] 值为 1, 否则为 0。

⊖ 有一些不重要的例外。例如, 一个过程可能有两个非嵌套的 **for** 语句, 都使用索引变量 i 。严格地讲, **for** 语句的索引变量的作用域是该 **for** 语句本身, 因此 i 其实是两个不同的变量。

作。在这种情况下,可假定共用数据块中的变量是全局的,在操作时过程使用的临时存储可假定是局部于过程的变量。对变量的作用域可能产生冲突的情况,本书将给出明确的说明。

习题

- 1.1 如果 a) 对于 $\varepsilon > 0$ 和有限例外的所有 n , 有 $f(n) \geq \varepsilon$
 b) 存在常数 $c_1 > 0$, $c_2 > 0$, 对于几乎所有的 $n \geq 0$, 有 $g(n) \leq c_1 f(n) + c_2$ 证明 $g(n)$ 是 $O(f(n))$ 的。
- 1.2 如果存在正的常数 c , 对所有的 n , 有 $f(n) \leq cg(n)$, 则记 $f(n) \leq g(n)$ 。证明若 $f_1 \leq g_1$, $f_2 \leq g_2$, 则 $f_1 + f_2 \leq g_1 + g_2$ 。关系 \leq 还有什么其他的性质吗?
- 1.3 编写完成下述任务的 RAM、RASP 和简化 ALGOL 语言程序:
 a) 给定输入 n , 计算 $n!$;
 b) 读入以结束符(0)结尾的 n 个正整数, 打印这 n 个整数的排序结果;
 c) 接受所有型为 1^*2^*0 的输入。
- 1.4 对(a)统一代价和(b)对数代价, 分析你所给出的针对习题 1.3 的答案的时间复杂度和空间复杂度。说明对输入“规模”的衡量。
- *1.5 编写计算 n^n 的统一代价时间复杂度为 $O(\log n)$ 的 RAM 程序并证明其正确性。
- *1.6 证明对数代价时间复杂度为 $T(n)$ 的 RAM 程序, 都存在一个等价的时间复杂度为 $O(T^2(n))$ 的无 MULT 或 DIV 指令的 RAM 程序。[提示: 使用偶数编号的寄存器作为工作存储器的子例程来仿真 MULT 和 DIV 指令。对于 MULT 指令, 说明: 如果 i 乘以 j , 则可以计算每个 $l(j)$ 部分积, 然后在 $O(l(j))$ 步内将它们加起来, 其中每步需要 $O(l(i))$ 时间。]
- *1.7 如果将 MULT 和 ADD 从指令集中移去, RAM 或 RASP 的计算能力会发生什么变化? 对于计算代价的影响又如何?
- **1.8 说明任何可被 RAM 接受的语言也可被无间接寻址指令的 RAM 所接受。[提示: 说明一整条图灵机带可以编码为单一的整数。由此, 任意的 TM 都可以使用有限个 RAM 的寄存器进行仿真。]
- 1.9 证明在一个常数因子范围内, RAM 和 RASP 对于(a)统一代价和(b)对数代价, 空间复杂度是等价的。
- 1.10 给定 9 个标量元素作为输入, 编写计算一个 3×3 行列式的直线状程序。
- 1.11 编写一个位操作的序列以计算两个 2 位整数的乘积。
- 1.12 说明直线状模型下任何由二进制输入和布尔运算符组成的 n 个语句所计算的函数集合皆可由包含 n 个布尔电路元件的组合逻辑电路实现。
- 1.13 说明任何布尔函数皆可由直线状程序计算。
- *1.14 假定用位向量 v_i 的集合来表示一个有 n 个顶点的图, 当且仅当从顶点 i 到顶点 j 有一条边, v_i 的第 j 个分量为 1。给出一个 $O_{bv}(n)$ 算法, 其计算结果是向量 p_1 , 当且仅当顶点 1 到顶点 j 之间存在一条路径, 其第 j 个分量的值为 1。可以使用的操作是对位向量进行运算的按位逻辑指令, 算术指令(对具有“整数类型”的变量), 将向量的某位设置为 0 或 1 的指令, 和一条特殊的指令(其含义为: 如果向量 v 的最左边的 1 位于第 j 个位置, 则将 j 赋值给整型变量 a , 如果 v 的所有分量都是 0, 则令 $a=0$)。
- *1.15 给出一台图灵机, 当给定的两个二进制整数分别位于带 1 和带 2 时, 在带 3 上打印它们的和。可以假定带的最左端放置了特殊符号#。
- 1.16 对于如图 1-21 所示的图灵机, 当输入分别为 a) 0010 和 b) 01110 时, 给出格局的变化序列。

- *1.17 给出一台完成如下工作的图灵机：
a) 当带 1 上以 0^n 开始时，在带 2 上打印 0^n 。
b) 接受型为 0^*10^n 的输入。
- 1.18 给出 TM 的状态集和下一移动函数，仿定理 1.3 的证明中的 RAM 指令 LOAD3。
- 1.19 给出时间复杂度为 $O(n)$ 的 RAM 算法，对于 n ，计算 2^{2^n} 。程序的 (a) 统一代价和 (b) 对数代价为何？
- *1.20 给定 $g(0, n) = n$ 和 $g(m, n) = 2^{g(m-1, n)}$ ， $m > 0$ ，定义函数 $g(m, n)$ 。给出对于给定的 n ，计算 $g(n, n)$ 的 RAM 程序。比较其统一代价和对数代价复杂度。
- 1.21 使用实参 i 和 $A[i]$ 执行 1.8 节的过程 INTERCHANGE，分别按名调用、按引用调用，结果一样吗？

研究性问题

- 1.22 图灵机是否能对定理 1.3 所阐述的在 $O(T^2(n))$ 的时间上界内模拟 RAM 计算的结果进行改进？

文献及注释

Shepherdson 和 Sturgis[1963]，Elgot 和 Robinson[1964] 以及 Hartmanis[1971] 都对 RAM 和 RASP 模型进行形式化的论述。本书给出的与 RAM 和 RASP 有关的大多数结果来自 Cook 和 Reckhow[1973] 的专利。

图灵机理论可参见 Turing[1936]。更多概念细节可参阅 Minsky[1967] 或 Hopcroft 与 Ullman[1969]，也是习题 1.8 的答案参考。图灵机的时间复杂度最早的研究文献有 Hartmanis 和 Stearns[1965]，关于空间复杂度的有 Hartmanis、Lewis 和 Stearns[1965] 以及 Lewis、Stearns 和 Hartmanis[1965]。计算复杂性的抽象研究开始于 Blum[1967]。综述文献可参见 Hartmanis 和 Hopcroft[1971] 以及 Borodin[1973a]。

Rabin[1972] 给出了计算的决策树模型的有趣扩充。

第2章 有效算法的设计

本章目的有两个。首先,介绍一些基本的、对许多问题有效算法的设计是非常有用的数据结构。其次,介绍一些在许多有效算法中很常见的“程序设计”技术,例如递归和动态规划等。

第1章考虑的是基本计算模型,虽然本书基于的模型是RAM,但一般不会根据其基本设备描述算法。因此,还介绍了简化ALGOL语言(1.8节)。然而这种语言较简单,还需要引入比数组更复杂的数据结构。本章首先介绍一些基本的数据结构,如读者熟悉的在有效的算法中经常用到的表和堆栈等。本章将讨论如何使用这些结构来表示集合、图和树。处理表的过程很简单,不熟悉的读者可查阅本章后面的基本参考资料,或者留意一下课后习题。

本章也包括递归的介绍。关于递归,很重要的一点是它能够使算法的描述简化。虽然本章的例子很简单,不能充分说明这一点,但是通过使用递归可减少细节的考虑,对简明扼要地阐述后面章节涉及的一些更复杂的算法是很有用的。递归本身并不一定使算法更有效,然而,当它与其他技术如平衡、分治、代数简化等结合在一起的时候,就会看到它常会使算法既高效又优美。

2.1 数据结构:表、队列和堆栈

假设读者已经熟知数学的基本概念如集合、关系以及基本数据类型如整数、字符串和数组等。这一小节将快速复习基本的表操作。

从数学上讲,表是从某个与应用相关的集合中所抽取的项的有限序列。算法描述经常包括一个表,并对表中的项进行添加或者删除操作。尤其是在表中的某个位置添加或者删除一项。基于这个原因,一般希望设计一种表的数据结构,在其中能够方便地任意增加或者删除项。

考虑表

$$\text{Item1, Item2, Item3, Item4} \quad (2-1)$$

这个表最简单的实现就是使用如图2-1所示的单链结构。结构中每一个元素包括两个存储单元。第一个存储单元包含项自身[⊖],第二个存储单元包含指向下一个元素的指针。一个可能的实现是使用如图2-2中称为NAME和NEXT的两个数组[⊖]。如果ELEMENT是数组的索引,那么NAME[ELEMENT]就是存储的项。如果有下一个表项的话,NEXT[ELEMENT]就是表中下一个项的索引。如果,ELEMENT是表中最后一个表项的索引,那么NEXT[ELEMENT]=0。

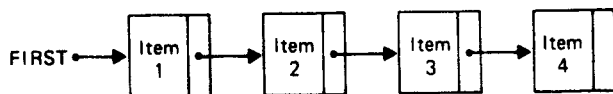


图2-1 链表

⊖ 如果项自身是一个复杂的结构,那么第一个存储单元可能包含一个指向该结构的指针。

⊖ 一个可选的(或等价的)观点是每一个元素有一个“单元”。每个单元有一个“地址”,该地址存储在为该元素所保留的一组寄存器中的第一个(可能是仅有的)寄存器中。在每个单元中有一个或多个“域”。这里域是NAME和NEXT,NAME[ELEMENT]和NEXT[ELEMENT]用来表示单元中域的内容,其中单元的地址是ELEMENT。

图 2-2 用 NEXT[0] 作为总是指向表的第一个元素的指针。注意, 项在数组 NAME 中的顺序与它们在表中的顺序并不相同。然而, 图 2-2 是图 2-1 的忠实的表示, 因为 NEXT 数组序列将项按它们在表(2-1)中的顺序排列。

接下来要考虑的过程是将一个元素插入一个表中。假设 FREE 是数组 NAME 和 NEXT 中未用单元的索引, POSITION 是表中某元素的索引, ITEM 要插入到表中该元素的后面。

```

procedure INSERT(ITEM, FREE, POSITION):
begin
    NAME[FREE] ← ITEM;
    NEXT[FREE] ← NEXT[POSITION];
    NEXT[POSITION] ← FREE
end

```

任何合理的到 RAM 代码的翻译会将 INSERT 操作的执行代价的结果时间, 该执行时间与表的大小无关。

例 2.1 假设希望在表(2-1)的 Item2 后面插入一项 Newitem 以创建表

Item1, Item2, Newitem, Item3, Item4

如果图 2-2 中每个数组的存储单元 5 都为空, 可以通过调用 INSERT(Newitem, 5, 3) 在 Item2 (位置是 3) 的后面插入 Newitem。INSERT 的三个语句将 NAME[5] 赋值为 Newitem, NEXT[5] 赋值为 4, NEXT[3] 赋值为 5, 如此建立如图 2-3 所示的数组。 □

	NAME	NEXT
0	—	1
1	Item 1	3
2	Item 4	0
3	Item 2	4
4	Item 3	2

图 2-2 有四个项的表的表示

	NAME	NEXT
0	—	1
1	Item 1	3
2	Item 4	0
3	Item 2	5
4	Item 3	2
5	Newitem	4

图 2-3 插入 Newitem 后的表

为了删除在存储单元 I 后的元素, 可设置 NEXT[I] 等于 NEXT[NEXT[I]]。如果删除成功, 删除元素的索引将放置在未用的存储单元表中。

同一个数组常会存储几个表。一般, 其中一个表是未用存储单元的列表, 可把它作为空闲列表。为了添加一项到表 A , 可修改过程 INSERT, 通过删除空闲列表的第一个单元得到一个未用单元。从表 A 删除一项, 可将相应的单元返回至空闲列表备用。

这种存储管理方法并不是唯一可用的方法, 但是, 一旦决定从哪里添加或者删除项, 这种方法就可以在有限次的操作内完成对表实施添加和删除项的操作。

其他关于表的基本操作是归并两个表成为一个表, 以及它的逆操作, 即在某个元素后面将其断开成为两个表。归并操作可通过增加另一个指向表的表指针使操作在有限时间内执行。这个指针给出了表中最后一个元素的索引, 从而避免为了找到最后一个元素而搜索整个表。如果能够在分割之前给出元素的索引, 分割操作的时间可以是有界的。

通过增加一个数组 PREVIOUS 能够在两个方向对表进行遍历。PREVIOUS[I] 的值是表中位置为 I 的项的前趋邻接项的位置。称具有这种性质的表为双向链接。对于双向链表, 不用给出前一项的位置, 就能够删除或者插入一项。

有时对表实施的操作是非常有限的。例如,可能仅在表尾添加或者删除项。也就是说,以后进先出的方式插入或者删除项。在这种情况下,一般可称表为堆栈或者压栈存储。

堆栈经常以一维数组实现。例如,表

Item1, Item2, Item3

可存储在如图 2-4 所示的数组 NAME 中。TOP 变量是一个指针,指向最后一个压入堆栈的项。为了给堆栈添加(PUSH)新的项,可让 TOP 为 TOP + 1,然后给 NAME[TOP]赋一个新值(因为数组 NAME 的长度为有限的 l ,在插入新值前,需要确定 $TOP < l - 1$)。要从堆栈的顶部删除(POP)一项,设置 TOP 为 TOP - 1 即可。注意,不需要将项从堆栈物理性删除掉。可以通过检查 TOP 的值是否小于零来检测堆栈是否为空。显然,执行 PUSH 和 POP 操作、测试堆栈是否为空的操作时间是与堆栈的元素数目无关的。

另一个特别的表形式是队列,此时表中的项是从一端(队头)添加而从另一端删除。和堆栈一样,一般使用一维数组构造一个队列。图 2-5 表示包含表中项 P、Q、R、S、T 的队列。两个指针分别指向队头和队尾。为了向队列添加(ENQUEUE)一个新项,和堆栈一样,可将 FRONT 设为 FRONT + 1,再将新项存储在 NAME[FRONT]。从队列中移除(DEQUEUE)一项时,需要将 REAR 设置为 REAR + 1。注意,这种技术使项的存取以先进先出的方式进行。

因为数组 NAME 的长度是有限的即 l ,指针 FRONT 和 REAR 最终将到达数组的端点。如果队列所表示的表的长度不会超过 l ,可以将 NAME[0] 处理为 NAME[$l - 1$]-后的元素。

存储在表中的项本身可以是复杂的结构。例如,在操作一个数组表时,不必实际添加或者删除数组,因为每次添加或者删除操作都需要与数组大小成比例的时间。一般添加或者删除指向数组的指针即可。由此可见,一般可以在固定的时间内添加或者删除一个复杂的结构,所花费的时间与这个结构的大小无关。

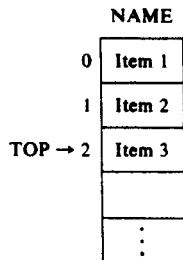


图 2-4 堆栈的实现

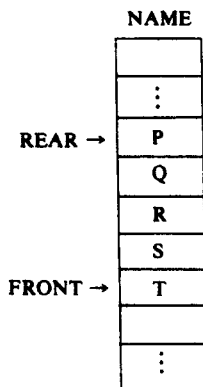


图 2-5 队列的一维数组实现

2.2 集合的表示

表通常用来表示集合。在这种表示法下,表示集合所需的内存规模与集合中元素的个数成比例。执行对集合的操作所需要的时间依赖于操作的特性。例如,假设 A 和 B 是两个集合,操作 $A \cap B$ 需要至少与两个集合的规模总和成比例的时间,原因是必须分别扫描表示 A 和 B 的集合^①。

同样,操作 $A \cup B$ 也至少需要与集合大小总和成比例的时间,原因是必须查找出现在两个集合中的相同元素,并且删掉其中之一。然而,如果 A 和 B 没有交集,可简单地合并两个表示 A 和 B 的表,此时 $A \cup B$ 操作的时间是与 A 和 B 的大小无关的。如果需要一种快速的方法确定给定的元素是否在给定集合中出现,不相交的集合的合并问题将会变得更加复杂。4.6 和 4.7 节将更全面地讨论这个问题。

除了表,另一种表示集合的方法是使用位向量。假设所讨论的域为 U (所有的集合都是它的子集),它有 n 个成员。先线性排列 U 中的元素。子集 $S \subseteq U$ 可表示为一个 n 位的向量 v_s ,当且仅当 U 中第 i 个元素是 S 的元素时, v_s 中的第 i 位为 1。一般称 v_s 是 S 的特征向量。

① 如果两个表有序,那么存在一个找到它们的交的线性时间算法。

位向量表示法具有在确定 U 的第 i 个元素是否属于一个集合的操作时间与集合大小无关的优点。而且, 集合的基本操作如并和交能够通过位向量的 \vee 和 \wedge 操作来执行。

如果不认为位向量操作是基本的(需要一个单位时间)运算, 那么可通过定义一个数组 A 来获得特征向量的效果, 当且仅当 U 的第 i 个元素在集合 S 中时, $A[i] = 1$ 。在这种表示方法下, 确定一个元素是否一个集合的成员是很简单的。它的缺点是交和并操作需要的时间是与 $|U|^\ominus$ 成比例的, 而不是与操作所涉及的集合大小成比例。类似地, 存储集合 S 的空间也与 $|U|$, 而不是与 $|S|$ 成比例。

2.3 图

现在介绍图的数学表示和常用的表示图的数据结构。

定义 一个图 $G = (V, E)$ 包括一个有限非空的顶点集 V , 和一组边的集合 E 。如果边是顶点的有序对 (v, w) , 那么称这个图为有向的, v 称为边 (v, w) 的尾, w 为边 (v, w) 的头。如果边是由不同的无序顶点对(集合)组成, 也表示为 (v, w) , 此时称图为无向的。[⊖]

在一个有向图 $G = (V, E)$ 中, 如果 (v, w) 是 E 中的一条边, 那么顶点 w 与顶点 v 邻接。边 (v, w) 从 v 到 w 。称 v 的邻接点的数量为 v 的(出)度。

在一个无向图 $G = (V, E)$ 中, 如果 (v, w) 是 E 中的一条边, 则 $(w, v) = (v, w)$, 即 (w, v) 和 (v, w) 是同一条边。如果 (v, w) [因此 (w, v) 也] 属于 E , 则说 w 与 v 邻接。顶点的度就是与它邻接的顶点数量。

有向图或无向图中的路径是边 $(v_1, v_2), (v_2, v_3), \dots, (v_{n-1}, v_n)$ 的序列。从 v_1 到 v_n 的路径长度为 $n - 1$ 。通常, 路径也可用路径上的顶点序列 v_1, v_2, \dots, v_n 表示。作为一个特例, 单个顶点的路径即从顶点到自身的路径, 长度为 0。如果除了第一个和最后一个顶点可能相同外, 路径中所有的边和顶点都不相同, 则路径称为是简单路径。环路是一个路径长度至少为 1 的简单路径, 它的起始和终止都在同一个顶点。注意, 在无向图中, 环路的长度至少为 3。

有几种常用的表示图 $G = (V, E)$ 的方法。其中之一是邻接矩阵, 即 $|V| \times |V|$ 的 0 和 1 矩阵 A 。在这个矩阵中, 当且仅当从顶点 i 到顶点 j 有一条边时, A 的第 ij 个元素 $A[i, j]$ 为 1。在经常需要查询某些边是否存在的图算法中, 图的邻接矩阵表示法是很方便的, 因为需要确定一条边是否存在的时间与 $|V|$ 和 $|E|$ 的大小无关, 是固定的。使用邻接矩阵的主要缺点是, 即使图只有 $O(|V|)$ 条边, 仍需要 $|V|^2$ 的存储空间。而简单直接地对邻接矩阵进行初始化的操作也需要 $O(|V|^2)$ 时间, 这使得对只有 $O(|V|)$ 条边的图进行操作时, 算法的时间复杂度仍需 $O(|V|)$ 。虽然存在解决这个问题的方法, 但总会有其他问题出现(见练习题 2.12), 使得基于邻接矩阵的时间复杂度为 $O(|V|)$ 的算法很少。

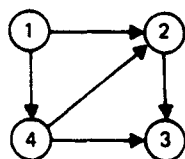
一个有趣的选择是用位向量来表示邻接矩阵的行和/或列。这种表示法可以使图算法有相当的效率。

另一个可能的图的表示法是通过表的方式。顶点 v 的邻接表是一个包括所有与 v 邻接的顶点 w 的表。一个图能用 $|V|$ 个邻接表表示, 每个顶点用一个表。

例 2.2 图 2-6a 表示有 4 个顶点的有向图。图 2-6b 是相应的邻接矩阵。图 2-6c 表示 4 个邻接表, 每个顶点一个表。例如, 由于存在着从顶点 1 到顶点 2 和顶点 4 的边, 因此顶点 1 的邻接表中有元素 2 和 4, 它们以图 2-1 的形式链接到一起。

⊖ 本书用 $|X|$ 表示集合 X (大小或容量) 中元素的数量。

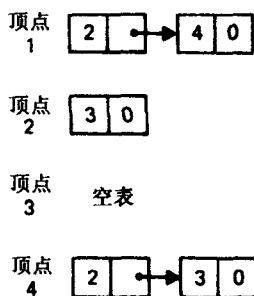
⊖ 注意, (a, a) 可能是有向图而不是无向图的边。



a) 有向图

	1	2	3	4
1	0	1	0	1
2	0	0	1	0
3	0	0	0	0
4	0	1	1	0

b) 邻接矩阵



c) 邻接表

	HEAD	NEXT
顶点 { 1		5
2		7
3		0
4		8
边 { 5	2	6
6	4	0
7	3	0
8	2	9
9	3	0

d) 邻接表的表格表示法

图 2-6 有向图和它的表示方法

图 2-6d 是邻接表的表格表示法。NEXT 数组的前 4 个存储单元存储着指向邻接表第一个顶点的指针，即 NEXT[i] 指向顶点 i 的邻接表的第一个顶点。注意，由于没有顶点在顶点 3 的邻接表上，所以 NEXT[3] = 0。数组 NEXT 的剩下部分表示图的边。数组 HEAD 包含邻接表中的顶点。这样，由于 NEXT[1] = 5，顶点 1 的邻接表从存储单元 5 开始。HEAD[5] = 2 表示有一条边 (1, 2)。NEXT[5] = 6 和 HEAD[6] = 4 表示边 (1, 4)。NEXT[6] = 0 表示没有其他以 1 为尾的边了。□

注意，用邻接表表示图时，需要的存储空间与 $|V| + |E|$ 成比例。当 $|E| \ll |V|^2$ 时，常用邻接表表示法。

如果图是无向的，那么每条边 (v, w) 都会表示两次，一次是在 v 的邻接表中，另一次是在 w 的邻接表中。在这种情况下，可能会增加一个 LINK 数组使无向图的副本间相互关联。这样，如果 i 是顶点 w 在顶点 v 的邻接表中的位置，LINK[i] 就是顶点 v 在顶点 w 的邻接表中的位置。

要方便地从无向图中删除边，可使用双向链接邻接表。这通常是必需的，因为即使删除边 (v, w) 是顶点 v 的邻接表上的第一条边，反方向的边可能会出现在顶点 w 的邻接表的中间。为了快速地从顶点 w 的邻接表中删除边 (v, w) ，必须能够迅速地在这个邻接表中找到它的前一条边的位置。

2.4 树

接下来介绍一种非常重要的有向图——树，还要考虑适当表示它的数据结构。

定义 称没有环路的有向图为有向无环图。(有向)树(有时称为根树)是一个满足下列性质的有向无环图：

1. 恰好只有一个顶点没有边进入, 该顶点称为根;
2. 除了根, 每个点都恰好有一入边;
3. 根到每个顶点都有一条路径(显然这条路径是唯一的)。

一个有向图若是一些树的集合, 则称为森林。森林和树是有向无环图的特殊情况, 它们出现是非常频繁的。下面是一些相关的术语。

定义 设 $F=(V, E)$ 是森林的图。如果 (v, w) 在 E 中, 则称 v 为 w 的父亲, w 是 v 的儿子。如果有一条从 v 到 w 的路径, 那么称 v 是 w 的祖先, w 是 v 的后代。若 $v \neq w$, 那么称 v 是 w 的严格祖先, w 是 v 的严格后代。一个没有严格后代的顶点叫做树叶。顶点 v 和它所有的后代叫做 F 的子树。顶点 v 叫做子树的根。

树中顶点 v 的深度是从根到 v 的路径长度。树中顶点 v 的高度是从 v 到树叶的最长路径的长度。树的高度是根的高度。顶点 v 在树中的层数是指树的高度减去 v 的深度。例如, 在图 2-7a 中, 顶点 3 的深度是 2, 高度为 0, 层数为 1。

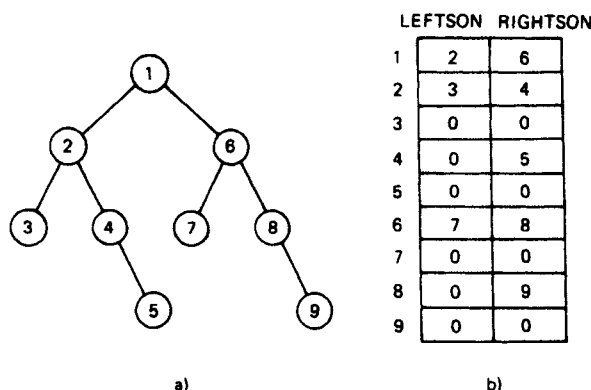


图 2-7 二叉树和它的表示方法

有序树是一棵顶点的儿子间有序的树。画有序树时, 假设每一个顶点的儿子是由左向右排序的。二叉树是如下定义的有序树:

1. 顶点的每个儿子都明确是左儿子或者右儿子;
2. 没有顶点有多于一个左儿子或多于一个右儿子。

以顶点 v 的左儿子为根的子树 T_l (如果存在) 称为 v 的左子树。类似, 以顶点 v 的右儿子为根的子树 T_r (如果存在) 称为 v 的右子树。 T_l 中的所有顶点均在 T_r 中所有顶点的左边。

二叉树常用两个数组 LEFTSON 和 RIGHTSON 表示。若二叉树的顶点用从 1 到 n 的整数表示, 那么, 当且仅当 j 是 i 的左儿子时, $\text{LEFTSON}[i] = j$ 。如果 i 没有左儿子, 那么 $\text{LEFTSON}[i] = 0$ 。 $\text{RIGHTSON}[i]$ 可以类似定义。

例 2.3 图 2-7a、b 给出一棵二叉树和它的表示。 □

定义 对于某个整数 k , 若每个深度均小于 k 的顶点既有左儿子又有右儿子, 而每个深度为 k 的顶点都是树叶, 则称这个二叉树是完全的。高度为 k 的完全二叉树恰有 $2^{k+1} - 1$ 个顶点。

一个高度为 k 的完全二叉树通常可用单个数组表示。数组的第一个位置包含树根。位置 i 的顶点的左儿子位于位置 $2i$, 右儿子位于位置 $2i + 1$ 。假如一个顶点在位置 $i > 1$, 则它的父亲的位置在 $\lfloor i/2 \rfloor$ 。

许多涉及树的算法经常以某种顺序遍历树(访问树的每个顶点)。存在着几种系统的遍历方式。三种常用的是前序、后序和中序。

定义 设 T 为一棵树, 有根 r 及其儿子 $v_1, v_2, \dots, v_k, k \geq 0$ 。在 $k=0$ 的情况下, 树只有一个顶点 r 。

T 的前序遍历如下递归地定义:

1. 访问根 r 。
2. 以前序依次访问根为 v_1, v_2, \dots, v_k 的子树。

T 的后序遍历如下递归地定义:

1. 以后序依次访问根为 v_1, v_2, \dots, v_k 的子树。
2. 访问根 r 。

对于二叉树, 中序遍历如下递归地定义:

1. 以中序访问根 r 的左子树(如果存在的话)。
2. 访问 r 。
3. 以中序访问 r 的右子树(如果存在的话)。

例 2.4 图 2-8 表示一棵二叉树, 顶点以前序(图 2-8a)、后序(图 2-8b)和中序(图 2-8c)的顺序进行标记。□

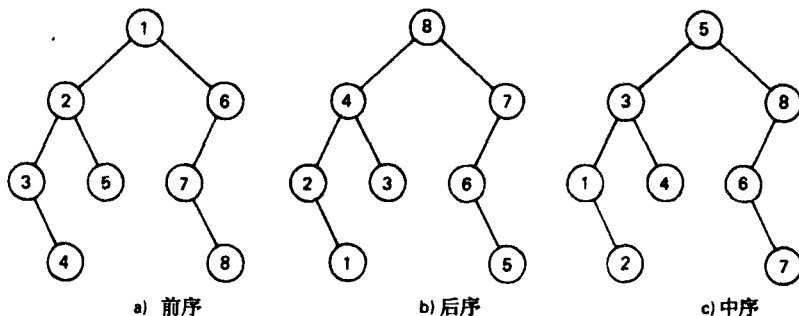


图 2-8 树的遍历

一旦通过遍历给顶点指派了标号, 那么, 通过标号查到顶点是很方便的。因此, 也可用顶点的标号 v 来表示该顶点。若以访问顺序来标号顶点, 那么标号就有一些有趣的性质。

在前序中, 以 r 为根的子树的所有顶点的标号不会比 r 小。更准确地说, 如果 D_r 是 r 的后代集合, 那么, 当且仅当 $r \leq v < r + |D_r|$ 时, v 在 D_r 中。通过将每个顶点 v 与前序标号和后代的数目联系起来, 可以很容易地确定顶点 w 是否 v 的后代。在初始指派前序标号并计算每个顶点的后代的数目后, w 是否 v 的后代的问题能够在固定时间内回答, 与树的大小无关。后序标号也有类似的性质。

二叉树的中序标号的性质是顶点 v 的左子树的每个顶点的标号比 v 小, 而且 v 的右子树的每个顶点的标号比 v 大。因此, 为了找到顶点 w , 将 w 与根 r 进行比较。如果 $w = r$, 那么 w 已经找到。如果 $w < r$, 在左子树重复这个过程; 如果 $w > r$, 在右子树重复这个过程。最终将找到 w 。遍历的这个特性将用在后面的章节中。

下面是关于树的最终定义。

定义 无向树是一个连通(任意两个顶点之间存在通路)、无环的无向图。有根无向树是显式指派一个顶点为根的无向树。

简单地使有向树的所有边变为无向, 就可以使有向树转化为一棵无向有根树。对于无向有根树, 可使用与有向树相同的术语和表示习惯。基本的数学上的差别是: 在有向树中所有的路径都是从祖先到后代, 而在无向有根树中, 路径是双方向的。

2.5 递归

一个直接或间接调用自身的过程叫做递归。跟非递归相比, 使用递归可以尽可能简明地描述算法。这一节将给出递归算法的例子, 并勾勒出递归是怎样在 RAM 上执行的。

考虑 2.4 节中给出的二叉树的中序遍历。我们希望在创建一个为二叉树顶点分配中序标号的算法能反映中序遍历的定义。下面给出这样的算法。注意, 这个算法通过递归调用自身来标记子树。

算法 2.1 二叉树顶点的中序标号。

输入: 用数组 LEFTSON 和 RIGHTSON 表示的二叉树。

输出: 称为 NUMBER 的数组, $NUMBER[i]$ 是顶点 i 的中序标号。

方法: 除了 LEFTSON, RIGHTSON 和 NUMBER, 算法还使用一个全局变量 COUNT, 它包含要分配给各顶点的中序标号。COUNT 初值为 1。参数 VERTEX 最初为根。递归调用图 2-9 所示的过程。

算法本身是:

```

procedure INORDER(VERTEX):
begin
1.  if LEFTSON[VERTEX]  $\neq$  0 then
    INORDER(LEFTSON[VERTEX]);
2.  NUMBER[VERTEX]  $\leftarrow$  COUNT;
3.  COUNT  $\leftarrow$  COUNT + 1;
4.  if RIGHTSON[VERTEX]  $\neq$  0 then
    INORDER(RIGHTSON[VERTEX])
end

```

图 2-9 中序遍历的递归过程

```

begin
    COUNT  $\leftarrow$  1;
    INORDER(ROOT)
end

```

□

使用递归有几个好处。首先, 通常递归程序更容易理解。如果上面的算法没有用递归, 则可能需要构建一个显示的机制来遍历树。树中向下的路径是没有什么问题的, 但是为了有返回祖先的能力, 则需要在堆栈中存储祖先, 而且操作堆栈的语句将会使得算法复杂化。同一个算法的非递归版本看起来可能如下。

算法 2.2 算法 2.1 的非递归版本。

输入: 同算法 2.1。

输出: 同算法 2.1。

方法: 在堆栈中存储从根到当前被搜寻的顶点的路径上所有未被标记的顶点的方式来遍历树。在从顶点 v 到 v 的左儿子的遍历过程中, 将 v 存储在堆栈中。在搜索完 v 的左子树后, 对 v 进行标记并将其从堆栈中弹出。接下来, 标记 v 的右子树。

在从顶点 v 到它的右儿子的过程中, 算法并没有将 v 放置在堆栈中, 这是因为, 标记完右子树后算法并不希望返回 v , 而是返回到 v 的还没标记的祖先(v 的最近的祖先 w , 此时 v 是 w 的左子树)。算法如图 2-10 所示。

□

```

begin
  COUNT ← 1;
  VERTEX ← ROOT;
  STACK ← empty;
left:  while LEFTSON[VERTEX] ≠ 0 do
        begin
          push VERTEX onto STACK;
          VERTEX ← LEFTSON[VERTEX];
        end;
center: NUMBER[VERTEX] ← COUNT;
        COUNT ← COUNT + 1;
        if RIGHTSON[VERTEX] ≠ 0 then
          begin
            VERTEX ← RIGHTSON[VERTEX];
            goto left;
          end;
        if STACK not empty then
          begin
            VERTEX ← top element of STACK;
            pop STACK;
            goto center;
          end;
        end;
end

```

图 2-10 非递归的中序算法

通过归纳二叉树顶点的标号，可容易地验证递归方法是正确的。类似地可证明非递归方法也是正确的，但是归纳假设并不像递归方法那样透明，需要额外关心堆栈操作以及二叉树的正确遍历。另一方面，递归可能付出的代价是对时间、空间复杂度的常数因子的影响。

很自然有递归算法怎样翻译成 RAM 代码的问题。根据定理 1.2，RAM 能够最多用慢常数因子的程序来模拟 RASP，因此讨论 RASP 的代码结构就可以了。这里要讨论一种比较直接的实现递归的技术。这种技术可以满足本书中用到的所有程序，但是它并没有覆盖所有可能出现的情况。

递归过程实现的核心是堆栈，堆栈中存储着每个没有终止的过程调用所使用的数据。也就是说，所有非全局数据都在堆栈中。可将堆栈划分为堆栈帧，它们是连贯的单元（寄存器）块。每个过程调用使用一个堆栈帧，这个堆栈帧的长度依赖于特定的调用过程。

假设正在执行过程 A。堆栈将如图 2-11 所示。如果 A 调用过程 B，将做如下操作：

1. 在堆栈的顶部放置一个适当大小的堆栈帧。帧的内部以 B 知道的方式排列：

- a) 指向调用 B 的实参的指针。[⊖]
- b) 清空用于存储 B 的局部变量的空间。
- c) 在结束对 B 的调用后，应执行例程 A 中 RASP 指令的地址（返回地址）。[⊖]

如果 B 是一个有返回值的函数，则指向 A 堆栈帧中存储函数返回值的单元的指针（值的地址）也放到 B 的堆栈帧中。

2. 控制权传给 B 的第一条指令。任何属于 B 的参数值或者局部标识符的地址都可以通过索引从 B 的堆栈帧中找到。

3. 当 B 结束，它通过以下步骤返回控制权给 A。

- a) 从栈顶获得返回地址。
- b) 如果 B 是一个函数，return 语句的表达式部分表示的值存储在堆栈中为存放值所规定的

⊖ 如果实参是一个表达式，则在 A 的堆栈帧中计算它，并把指向这个值的指针放置到 B 的帧中。如果实参是一个结构如数组，使用一个指向该结构的第一个字的指针就足够了。

⊖ 这是一个指向返回地址的跳转，这对 RAM 来说是难以实施的（虽然完全有可能），因此使用 RASP 模型。

单元中。

c) 过程 B 的堆栈帧从堆栈中弹出。这使得为过程 A 创建的堆栈帧位于堆栈的顶部。

d) A 从返回地址给出的单元处继续执行。

例 2.5 考虑算法 2.1 的过程 INORDER, 当它把 LEFTSON[VERTEX] 当作实参调用自己时, 它在堆栈中与 VERTEX 新值的地址共同存储的是标示调用结束的返回地址, 即以行 2 的语句继续执行程序。因此, 无论 VERTEX 在过程定义的任何地方出现, 变量 VERTEX 都能高效地替代为 LEFTSON[VERTEX]。

在某种意义上说, 上述操作模拟了非递归算法 2.2。然而, 认识到以 RIGHTSON[VERTEX] 为实参

的 INORDER 调用的完成也结束了调用过程的执行。因此, 当实参是 RIGHTSON[VERTEX] 的时候, 没有必要在堆栈中存储返回地址或者存储 VERTEX。□

过程调用所需时间与计算实参的时间以及在堆栈中存储指向值的指针花费成比例。返回语句所用的时间肯定不会超过这个时间。

考虑递归过程的集合花费的时间, 通常最简单的是计算一个调用过程的调用所花费的时间。这样作为输入尺度的函数, 可以限制每个过程调用花费的时间, 这不包括它调用的过程所花费的时间。可把所有过程调用的限制时间加起来, 给出一个总的代价的上限。

为了计算递归算法的时间复杂度, 可使用递推方程。作为输入参数 n 的函数, 函数 $T_i(n)$ 与第 i 个过程相联系, 表示第 i 个过程的执行时间。通常, 过程 i 所调用的过程的执行时间可表示为 $T_i(n)$ 的递推方程。然后, 求解递推方程组得到结果。若程序只涉及一个过程, 此时 $T(n)$ 依赖于 $T(m)$ 的值, 而 m 属于小于 n 的一个有限集合。下一节将讨论经常遇到的一些递归问题的解决方法。

记住, 在这里和任何其他地方, 所有的代价分析都假设使用统一代价函数。如果使用对数代价函数, 用来执行递归过程的堆栈长度可能影响时间复杂度分析。

2.6 分治法

一个解决问题的通常方法是将问题分为几个小的部分, 找到每部分的解决方案, 然后, 把每部分的解决方案结合起来, 组成整体的解决方案。尤其使用递归时, 如果子问题是原问题的较小的版本, 这种方法通常可以提供有效的解决方案。通过分析递推方程的结果, 下面两个例子说明了这种技术。

考虑在包含 n 个元素的集合 S 中找到最大元素和最小元素的问题。为了简单起见, 假设 n 是 2 的幂。找到最大和最小元素的一种显而易见的方式是分别找到它们。例如, 以下过程在 S 的元素中经过 $n-1$ 次比较, 寻找 S 中的最大元素。

```
begin
    MAX ← any element in S;
    for all other elements  $x$  in  $S$  do
        if  $x > \text{MAX}$  then MAX ←  $x$ 
end
```

类似地, 能够用 $n-2$ 次比较剩下的 $n-1$ 个元素, 以找到最小元素, 假设 $n \geq 2$, 通过总共 $2n-3$ 次比较, 找到最大和最小元素。

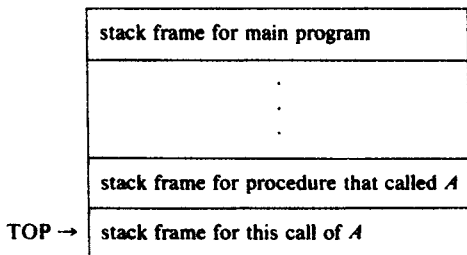


图 2-11 递归过程调用堆栈

```

    procedure MAXMIN(S):
1.  if ||S|| = 2 then
        begin
2.      let S = {a, b};
3.      return (MAX(a, b), MIN(a, b))
        end
    else
        begin
4.      将 S 分为两个子集 S1 和 S2, 每个子集都有一半的元素;
5.      (max1, min1) ← MAXMIN(S1);
6.      (max2, min2) ← MAXMIN(S2);
7.      return (MAX(max1, max2), MIN(min1, min2))
        end
    end

```

图 2-12 找到 MAX 和 MIN 的过程

分治方法将集合 S 分为两个子集 S_1 和 S_2 , 每个子集有 $n/2$ 个元素。通过算法的递归应用, 算法将会找到两个子集的最大和最小元素。可通过对 S_1 和 S_2 的最大和最小元素进行多于两次的比较而计算出 S 的最大和最小元素。详细算法如下:

算法 2.3 找到一个集合的最大和最小元素。

输入: 有 n 个元素的集合 S , n 是 2 的幂, 并且, $n \geq 2$ 。

输出: S 的最大和最小元素。

方法: 应用递归过程 MAXMIN 于集合 S 。MAXMIN 有一个变量[⊖]是集合 S , $|S| = 2^k$, $k \geq 1$, MAXMIN 返回一对值 (a, b) , a, b 分别是 S 中的最大、最小元素。过程 MAXMIN 在图 2-12 中给出。□

注意, 需要在 S 的元素间进行比较的步骤仅是步骤 3 和步骤 7, 步骤 3 中比较 S 的两个元素, 步骤 7 则必须比较 max1 和 max2 以及 min1 和 min2。设 $T(n)$ 是 MAXMIN 在 n 个元素的数据集 S 中找到最大和最小元素所需的比较次数。显然, $T(2) = 1$ 。如果 $n > 2$, $T(n)$ 是对大小为 $n/2$ 的集合的两次调用 MAXMIN 中所需的比较次数 (行 5、6) 加上第 7 行的两次比较的和, 即:

$$T(n) = \begin{cases} 1 & \text{对于 } n=2 \\ 2T(n/2) + 2 & \text{对于 } n>2 \end{cases} \quad (2-2)$$

函数 $T(n) = (3/2)n - 2$ 是递归式 (2-2) 的解。显然, $n=2$ 时式 (2-2) 满足, 如果 $n=m$ ($m \geq 2$) 时, 式 (2-2) 满足, 那么

$$T(2m) = 2\left(\frac{3m}{2} - 2\right) + 2 = \frac{3}{2}(2m) - 2$$

因此, $n=2m$ 时, 式 (2-2) 也满足。这样, 通过对 n 的归纳, 可以知道只要 n 是 2 的幂, $T(n) = (3/2)n - 2$ 可使式 (2-2) 满足。

可以证明在 n 个元素的集合 S 中寻找最大和最小元素, 至少需要进行 $(3/2)n - 2$ 次比较。这样, 当 n 是 2 的幂时, 算法 2.3 在 S 的元素间比较的次数是最优的。

在前面的例子中, 分治方法的使用减少了比较次数的常数因子倍。接下来的例子将说明如何通过利用分治技术实际减少算法的渐近增长率。

考虑两个 n 位数的乘法。传统的方法需要 $O(n^2)$ 次位操作。下面的方法大约需要 n^{\log_3} 或者它的近似值 $n^{1.59}$ 次位操作。[⊖]

⊖ 由于仅考虑比较, 参数的传递方法不那么重要。然而, 如果集合 S 用数组表示, 可通过传递指向 S 的子集的第一个和最后一个元素的指针有效地调用 MAXMIN, 其中 S 的子集是数组中连续的元素。

⊖ 除非另有声明, 本书所有对数均以 2 为底。

设 x 和 y 为两个 n 位数。同样, 为了简化起见, 假设 n 是 2 的幂。如图 2-13, 将 x 和 y 分为两半。如果将每一半为 $(n/2)$ 位数, 那么, 可如下表示乘积:

$$xy = (a2^{n/2} + b)(c2^{n/2} + d) \quad (2-3)$$

$$= ac2^n + (ad + bc)2^{n/2} + bd$$

方程(2-3)通过 4 个 $(n/2)$ 位数的乘法与一些加法和移位(通过 2 的幂的乘法)来计算 x 和 y 的乘积。 x 和 y 的乘积 z 也可以用以下程序表示。

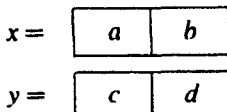


图 2-13 位串的切分

```

begin
    u ← (a + b) * (c + d);
    v ← a * c;
    w ← b * d;
    z ← v * 2^n + (u - v - w) * 2^{n/2} + w
end
    
```

(2-4)

先忽略进位的时候 $a + b$ 和 $c + d$ 可能是 $(n/2 + 1)$ 位数的情况, 假设它们仅是 $n/2$ 位数。为了使两个 n 位数相乘, 这个方案仅需要 $(n/2)$ 位数的 3 次乘法, 加上一些加法和移位。可以递归地使用乘法来计算乘积 u , v 和 w 。加法和移位需要时间为 $O_B(n)$ 。因此, 两个 n 位数相乘的时间复杂度的上界为

$$T(n) = \begin{cases} k & \text{对于 } n = 1 \\ 3T(n/2) + kn & \text{对于 } n > 1 \end{cases} \quad (2-5)$$

这里, k 是常数, 反映了式(2-4)中的加法和移位。递推式(2-5)的解的上界如下:

$$3kn^{\log_3} \approx 3kn^{1.59}$$

实际上式(2-5)的解是

$$T(n) = 3kn^{\log_3} - 2kn$$

因为 n 是 2 的幂, 通过对 n 进行归纳, 可以证明上式。归纳基础 $n = 1$ 是平凡的。如果 $n = m$ 时, $T(n) = 3kn^{\log_3} - 2kn$ 使式(2-5)满足, 那么对于递推步:

$$\begin{aligned}
 T(2m) &= 3T(m) + 2km \\
 &= 3[3km^{\log_3} - 2km] + 2km \\
 &= 3k(2m)^{\log_3} - 2k(2m)
 \end{aligned}$$

由此, $T(n) \leq 3kn^{\log_3}$ 。注意在归纳中, 若使用 $3kn^{\log_3}$ 而不是 $3kn^{\log_3} - 2kn$ 会导致失败。

为了完成乘法算法, 还必须考虑 $a + b$ 和 $c + d$ 是 $(n/2 + 1)$ 位数的情况。因此 $(a + b)(c + d)$ 的乘积不能直接用大小为 $n/2$ 的问题的递归算法计算。必须将 $a + b$ 写为 $a_12^{n/2} + b_1$, 其中 a_1 是 $(a + b)$ 的最高位, b_1 为剩下的位。类似地, 将 $c + d$ 写为 $c_12^{n/2} + d_1$ 。 $(a + b)(c + d)$ 的乘积可以表示为

$$a_1c_12^n + (a_1d_1 + b_1c_1)2^{n/2} + b_1d_1 \quad (2-6)$$

项 b_1d_1 可递归应用大小为 $n/2$ 的乘法算法进行计算。式(2-6)里面的其他乘法运算, 或者包括一位 a_1 和 c_1 , 或者是 2 的幂作为参数, 因此能够在 $O_B(n)$ 时间内完成。

例 2.6 上述快速整数乘法算法不仅可应用在二进制上, 也可以应用到十进制数上。下面的计算说明了这个方法。

$$\begin{array}{lll}
 x = 3141 & a = 31 & c = 59 \\
 y = 5927 & b = 41 & d = 27 \\
 a + b = 72 & c + d = 86 & \\
 u = (a + b)(c + d) = 72 \times 86 = 6192 & &
 \end{array}$$

$$\begin{aligned}
 v &= ac = 31 \times 59 = 1829 \\
 w &= bd = 41 \times 27 = 1107 \\
 xy &= 18290000^{\oplus} + (6192 - 1829 - 1107) \times 100 + 1107 \\
 &= 18616707
 \end{aligned}$$

□

注意, 基于式(2-4)的算法用三个加法和减法代替了一个乘法[与式(2-3)比较]。这个替换为什么能够导致渐近高效的直观原因是乘法比加法的执行难度更大, 对于足够大的 n , 不管用什么样的(已知)算法, 任何固定的 n 位的加法都比 n 位的乘法需要的时间少。初步看来, 似乎 $(n/2)$ 位数的乘的次数从 4 减到 3, 最好情况下, 能够减少总时间的 25%。然而, 由于式(2-4)可递归地应用于计算 $(n/2)$ 位、 $(n/4)$ 位、... 的乘积。由于每级都节省 25%, 因此渐近时间复杂度得到提高。

一个过程的时间复杂度是由于子问题的数量和大小决定的, 而将问题分解为子问题的总工作量也会对其产生次要影响。形如式(2-2)和式(2-5)的递归在递归分治算法分析中经常出现, 下面考虑在通常状况下的结果。

定理 2.1 设 a, b, c 是非负常数。递推式

$$T(n) = \begin{cases} b & \text{对于 } n=1 \\ aT(n/c) + bn & \text{对于 } n>1 \end{cases}$$

其中 n 是 c 的幂, 解是

$$T(n) = \begin{cases} O(n) & \text{若 } a < c \\ O(n \log n) & \text{若 } a = c \\ O(n^{\log_c a}) & \text{若 } a > c \end{cases}$$

证明: 如果 n 是 c 的幂, 那么

$$T(n) = bn \sum_{i=0}^{\log_c n} r^i, \text{ 其中 } r = a/c$$

如果 $a < c$, 那么 $\sum_{i=0}^{\infty} r^i$ 收敛, 因此, $T(n)$ 为 $O(n)$ 。如果 $a = c$, 那么和式中的每项都是 1, 并且有 $O(\log n)$ 项。因此, $T(n)$ 为 $O(n \log n)$ 。最后, 如果 $a > c$, 那么

$$bn \sum_{i=0}^{\log_c n} r^i = bn \frac{r^{1+\log_c n} - 1}{r - 1}$$

其为 $O(a^{\log_c n})$ 或者等于 $O(n^{\log_c a})$ 。□

从定理 2.1 可以看出分解一个问题(使用线性工作量)为一半大小的两个子问题, 将得到一个 $O(n \log n)$ 的算法。如果子问题的数量是 3、4 或者 8, 那么算法的阶分别为 n^{\log_3} 、 n^2 或者 n^3 。另一方面, 将问题分割为 4 个大小为 $n/4$ 的子问题, 也得到 $O(n \log n)$ 的算法, 而 9 和 16 个子问题使得算法的阶分别为 n^{\log_3} 和 n^2 。因此, 将整数分为 4 部分, 而且能够根据 8 个或者更少的乘法来表示整数的乘法, 可以得到整数乘法的渐近快速算法。也存在着一些重要的递推关系, 其分割问题的工作量与问题大小是不成比例的, 其中有些将在习题中出现。

在 n 不是 c 的幂的情况下, 通常能够将大小为 n 的问题嵌入到大小为 n' 的问题中, 其中 n' 大于等于 n 的最小 c 的幂。因此, 定理 2.1 的渐近增长率对于任意的 n 成立。实际上, 我们经常会设计将任意大小的问题分割为大小接近 c 的若干子问题的递归算法。这些算法通常比那些通过假设输入大小是 c 的某一幂而得到的算法更有效(相差一个常数因子)。

2.7 平衡

分治技术将问题分解为大小相等的若干子问题的两个例子并不是一个巧合。好的算法设计的一个基本原则是保持平衡。为了解释这个原则, 利用一个排序的例子来比较将问题分解为大小相等的

⊕ v 必须移动 4 个十进制位, $u-v-w$ 则移动两位。

子问题和大小不相等的子问题的效果差异。读者请不要从例子中推定平衡仅在分治技术中 useful。第 4 章的几个例子说明对子树规模进行平衡或者对操作的代价进行平衡将会使算法更加高效。

考虑 n 个整数按非递减次序进行排序的问题。也许, 最简单的排序方法是通过扫描序列定位最小的元素, 然后交换最小的元素与第一个元素的位置。这个过程在剩下的 $n-1$ 个元素间重复, 并将第二小的元素置于第二个位置。在剩下的 $n-2, n-3, \dots, 2$ 个元素排序中重复这个过程。

对于待排序元素间的比较次数, 由算法可得递归式

$$T(n) = \begin{cases} 0 & \text{对于 } n=1 \\ T(n-1) + n-1 & \text{对于 } n>1 \end{cases} \quad (2-7)$$

式(2-7)的解是 $T(n) = n(n-1)/2$, 即 $O(n^2)$ 。

虽然这个排序算法可看作分为不等的两部分的分治算法的递归应用, 但是, 对于大 n , 它不够有效。为了设计一个渐近有效的排序算法, 需要达到平衡。将问题分为两部分, 不是一部分为 1, 另一部分为 $n-1$ 的做法, 而是将问题分解为接近问题大小的一半的两个子问题。这可通过归并排序方法完成。

考虑一个整数序列 x_1, x_2, \dots, x_n 。再次为了简化问题, 假设 n 是 2 的幂。对该序列排序的一种方法是将其分解为两个序列 $x_1, x_2, \dots, x_{n/2}$ 和 $x_{(n/2)+1}, \dots, x_n$, 排序这两个序列, 然后归并。“归并”的意思是归并两个有序的序列为一个有序的序列。

算法 2.4 归并排序。

输入: 数值序列 x_1, x_2, \dots, x_n , 其中 n 是 2 的幂。

输出: 输入序列的排列 y_1, y_2, \dots, y_n , 满足 $y_1 \leq y_2 \leq \dots \leq y_n$ 。

方法: 利用过程 MERGE(S, T) 将两个有序序列 S 和 T 作为输入, 产生包含 S 和 T 的所有元素的有序序列作为输出序列。因为 S 和 T 自身有序, MERGE 最多需要 S 和 T 的长度之和减 1 次比较。通过反复对 S 和 T 中剩余元素中的最大元素进行比较, 选择其中较大的元素, 并删除所选元素。此过程一直持续, 直到 S 和 T 为空。

也可以利用图 2-14 的过程 SORT(i, j) 对序列 x_i, x_{i+1}, \dots, x_j 排序, 假设对 $k \geq 0$, 子序列长度为 2^k 。

为了排序给定的序列 x_1, x_2, \dots, x_n , 可调用 SORT($1, n$)。 □

在计算比较次数时, 算法 2.4 的递推关系为:

$$T(n) = \begin{cases} 0 & \text{对于 } n=1 \\ 2T(n/2) + n-1 & \text{对于 } n>1 \end{cases}$$

根据定理 2.1, 它的解为 $T(n) = O(n \log n)$ 。对于 n 很大的情况, 平衡子问题的大小带来的益处是很大的。类似的分析显示, 不仅是比较次数, 过程 SORT 所需的总时间也是 $O(n \log n)$ 。

```

procedure SORT( $i, j$ ):
  if  $i = j$  then return  $x_i$ 
  else
    begin
       $m \leftarrow (i + j - 1) / 2$ ;
      return MERGE(SORT( $i, m$ ), SORT( $m + 1, j$ ))
    end

```

图 2-14 归并排序

2.8 动态规划

如果通过合理的努力能够将问题分解为子问题, 而且子问题大小的总和能够保持很小, 递

归技术是很有用的。回忆定理2.1, 如果子问题大小的总和是 an , 常数 $a > 1$, 递归算法在时间复杂度上可能是多项式。然而, 如果划分大小为 n 的问题, 却导致 n 个大小为 $n-1$ 的问题, 那么算法的复杂度可能会指数增长。在这种情况下, 一种称为动态规划的表技术常会使算法更有效。

本质上, 动态规划计算了所有子问题的解。计算过程是从小的子问题到大一点的子问题, 并将结果存储在表中。这种方法的优点是一旦解决了子问题, 就将结果存储下来, 不需要重复计算。这个技术很容易用一个简单的例子解释。

考虑 n 个矩阵乘积的计算

$$M = M_1 \times M_2 \times \cdots \times M_n$$

这里每个 M_i 都是一个 r_{i-1} 行 r_i 列的矩阵。无论使用什么样的矩阵乘法算法, 矩阵相乘的顺序都会对计算 M 所需操作的总数有很大的影响。

例2.7 假设 $p \times q$ 矩阵和 $q \times r$ 矩阵相乘, 在“通常的”算法中, 需要 pqr 次操作, 考虑乘积

$$M = \begin{matrix} M_1 & \times & M_2 & \times & M_3 & \times & M_4 \\ [10 \times 20] & & [20 \times 50] & & [50 \times 1] & & [1 \times 100] \end{matrix} \quad (2-8)$$

这里每个 M_i 的维数标在括号中。以顺序

$$M = M_1 \times (M_2 \times (M_3 \times M_4))$$

计算 M , 需要 125 000 次操作, 然而, 以顺序

$$M = (M_1 \times (M_2 \times M_3)) \times M_4$$

计算 M , 仅需要 2200 次操作。□

为了最小化操作数量, 尝试所有可能的计算 n 个矩阵乘积的顺序是一个指数复杂度的过程 (见习题2.31), 这个过程当 n 很大的时候是不实用的。然而, 动态规划提供了一个 $O(n^3)$ 的算法。设 m_{ij} 是计算 $M_i \times M_{i+1} \times \cdots \times M_j$ 的最小代价, $1 \leq i \leq j \leq n$ 。显然,

$$m_{ij} = \begin{cases} 0 & \text{若 } i=j \\ \min_{i \leq k < j} (m_{ik} + m_{k+1,j} + r_{i-1}r_kr_j) & \text{若 } j > i \end{cases} \quad (2-9)$$

其中 m_{ik} 是计算 $M' = M_i \times M_{i+1} \times \cdots \times M_k$ 的最小代价。第2项 $m_{k+1,j}$ 是计算

$$M'' = M_{k+1} \times M_{k+2} \times \cdots \times M_j$$

的最小代价。而第3项是 M' 和 M'' 的乘积的代价。注意 M' 是一个 $r_{i-1} \times r_k$ 的矩阵, M'' 是 $r_k \times r_j$ 的矩阵。方程(2-9)说明考虑了 i 和 $j-1$ 之间所有可能的 k 值, m_{ij} 是这三项和的最小值, 其中 $j > i$ 。

动态规划方法以递增下标值的方式计算 m_{ij} 的值。首先对所有的 i 计算 m_{ii} , 然后对所有的 i 计算 $m_{i,i+1}$, 接下来是 $m_{i,i+2}$, 依此类推。这样, 当计算 m_{ij} 时, 式(2-9)中的 m_{ik} 和 $m_{k+1,j}$ 都能够用到。原因是, 如果 k 在 $i \leq k < j$ 范围内, $j-i$ 必须严格的比 $k-i$ 和 $j-(k+1)$ 大。下面给出了算法。

算法2.5 计算 n 个矩阵乘法 $M_1 \times M_2 \times \cdots \times M_n$ 的最小代价顺序的动态规划算法。

输入: r_0, r_1, \dots, r_n , 其中 r_{i-1} 和 r_i 是矩阵 M_i 的维数。

输出: M_i 相乘的最小代价, 假设矩阵 $p \times q$ 和矩阵 $q \times r$ 相乘需要 pqr 次操作。

方法: 算法如图2-15所示。□

例2.8 将算法应用到式(2-8)中的4个矩阵串中计算 m_{ij} 的值, 如图2-16所示, 其中 r_0, \dots, r_4 分别为 10, 20, 50, 1, 100。这样, 计算乘积所需的最小操作数目是 2200。乘法操作的顺序可通过记录确定, 表格(2-9)中的每个条目的 k 值给出了最小的操作数目。□

```

begin
1.   for i ← 1 until n do  $m_{ii} \leftarrow 0$ ;
2.   for l ← 1 until n - 1 do
3.     for i ← 1 until n - l do
4.       begin
5.          $j \leftarrow i + l$ ;
6.          $m_{ij} \leftarrow \text{MIN}_{i \leq k < j} (m_{ik} + m_{k+1,j} + r_{i-1} * r_k * r_j)$ 
7.       end;
8.   write  $m_{1n}$ 
9. end

```

图 2-15 矩阵乘法排序的动态规划算法

$m_{11} = 0$	$m_{22} = 0$	$m_{33} = 0$	$m_{44} = 0$
$m_{12} = 10\ 000$	$m_{23} = 1000$	$m_{34} = 5000$	
$m_{13} = 1200$	$m_{24} = 3000$		
$m_{14} = 2200$			

图 2-16 计算 $M_i \times M_{i+1} \times \cdots \times M_j$ 的代价

2.9 后记

本章涉及一些设计有效算法的基本技术。可以看到高级数据结构如表、队列和堆栈的使用会使算法设计者从琐碎的工作，如操纵指针等解放出来，而专注于算法全局结构的考虑。还可以看到递归和动态规划的功能之强大，通常使算法优美而自然。本章也给出了像分治和平衡这样一些普遍的原则。

当然这些并不是仅有的可用技术，但它们比较重要。本书的其他章节还将讨论其他技术。其范围将涉及问题的适当表示以及较好的操作顺序的选择。也许对一个好的程序设计者最重要的原则是永不满足。设计者应不断从不同的角度验证一个问题，直到确信他的算法是满足需求的最合适的算法为止。

习题

- 2.1 选择一个双向链表的实现。写出插入、删除一个数据项的简化 ALGOL 算法。确定你的程序在删除第一个和/或最后一个数据项以及表为空时可以正确执行。
- 2.2 编写一个颠倒表中数据项顺序的算法。并证明该算法能正确运行。
- 2.3 编写算法，实现 2.1 节中提到的操作 PUSH, POP, ENQUEUE 和 DEQUEUE。不要忘记检验指针是否到达堆栈或队列的数组尾部。
- 2.4 编写出检测队列是否为空的条件。假设 2.1 节中的数组 NAME 大小为 k ，那么队列可存储多少元素？画图说明当队列 (a) 为空，(b) 只有一个元素，(c) 队列满时，队列和指针 FRONT 及 REAR 的位置。
- 2.5 编写一个删除无向图中 v 的邻接表的第一条边 (v, w) 的算法。假设邻接表是双向链接的，如 2.3 节中描述的那样，可使用 LINK 在 w 的邻接表中定位 v 。
- 2.6 编写一个构造无向图的邻接表的算法。每条边 (v, w) 表示两次，一次是在 v 的邻接表中，

一次是在 w 的邻接表中。每条边的两个副本必须链接到一起，当删除一个，另一个也可以方便地删除。假设输入是边的表。

*2.7 (拓扑排序) 设 $G=(V, E)$ 是有向无环图。写一个为 G 的顶点赋整数值算法，如果存在一条从顶点 i 到顶点 j 的有向边，则 i 比 j 小。[提示：一个无环图肯定存在一个顶点没有入边。为什么？解决问题的一个方法是查找没有边进入的顶点，给这个顶点赋最小值，然后将它和它的所有出边从图中删除。在得到的图中重复这个过程，赋下一个最小值，依此类推。为了使算法有效，即时间复杂度为 $O(|E| + |V|)$ ，必须避免在每一个新产生的图中查找没有进入边的顶点。]

*2.8 设 $G=(V, E)$ 是有两个指定顶点，起始顶点和末尾顶点的有向无环图。写一个算法，找到一组从起始顶点到末尾顶点的路径，使得

1) 除了起始或末尾顶点，没有其他顶点在两条路径上出现；

2) 没有其他的满足条件(1)的路径可添加到集合中。

注意可能会有很多组路径满足上述条件。不需要找到路径最多的一组，只要满足上述条件就可以。算法执行时间必须为 $O(|E| + |V|)$ 。

*2.9 (稳定婚姻问题) 设 B 是 n 个男孩的集合， G 是 n 个女孩的集合。每个男孩将女孩从 1 到 n 排序，每个女孩将男孩从 1 到 n 排序。配对是男孩和女孩是一一对应的。如果每两个男孩 b_1 和 b_2 以及与他们配对的女孩 g_1 和 g_2 满足以下两个条件，则称配对是稳定的：

1) b_1 的排序中 g_1 比 g_2 高，或者 g_2 的排序中 b_2 比 b_1 高，

2) b_2 的排序中 g_2 比 g_1 高，或者 g_1 的排序中 b_1 比 b_2 高，

证明稳定配对总是存在的，写出一个算法找到一个这样的配对。

2.10 考虑每个顶点都有名字的二叉树。写一个算法按(a)前序、(b)后序和(c)中序，打印出名字。

*2.11 写一个算法计算带有 $+$ 和 \times 的算术表达式的(a)前缀波兰、(b)中缀和(c)后缀波兰表达式。

*2.12 开发一种技术，使得第一次存取矩阵项时，该项初始化为 0，从而消除初始化邻接矩阵的时间 $O(|V|^2)$ 。[提示：保持每个初始化过的项的指针为指向堆栈中后面的指针。每存取一项，通过确定那个项的指针指向堆栈的活动区域，而且栈的后面指针指向该项，从而验证该项的内容不是随机的。]

2.13 对图 2-17 的二叉树模拟算法 2.1 和算法 2.2。

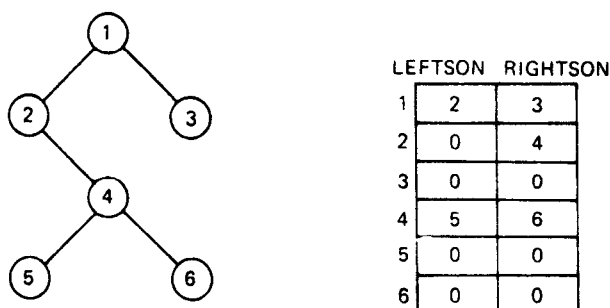


图 2-17 二叉树

*2.14 证明算法 2.2 是正确的。

- *2.15 (汉诺塔) 汉诺塔问题包括 3 个标桩 A、B、C 和 n 个不同大小的方盘。开始时, 方盘以递减的尺寸顺序堆积在标桩 A 上, 最大的方盘在底部。问题是每次将一个方盘从 A 移到 B, 并且没有方盘放置在比自身小的方盘上。标桩 C 能够用于暂时存放方盘。编写一个递归算法解决这个问题。以一个方盘移动所耗费的时间为单位, 算法的执行时间是多少?
- **2.16 求解习题 2.15 的非递归算法。考虑哪个算法更容易理解并证明是正确的?
- *2.17 证明 $2^n - 1$ 次移动对于解决习题 2.15 是充分和必要的。
- 2.18 写一个产生整数 1 到 n 的所有排列的算法。[提示: 整数 1 到 n 的排列的集合能够通过整数 1 到 $n-1$ 的每个排列中的可能位置插入 n 而得到。]
- 2.19 写一个求出二叉树高的算法。假设树如图 2-7b 所示。
- 2.20 写一个计算树中每个顶点的后代数量的算法。
- **2.21 考虑有两个输入和两个输出的双位置开关, 如图 2-18 所示。在一个位置输入 1 和 2 与输出 1 和 2 分别连接。在另一个位置输入 1 和 2 分别与输出的 2 和 1 连接。用这种开关, 设计一个有 n 个输入和 n 个输出的网络, 它能得到输入的 $n!$ 个任意可能的排列。网络必须用少于 $O(n \log n)$ 个开关。[提示: 利用分治方法。]

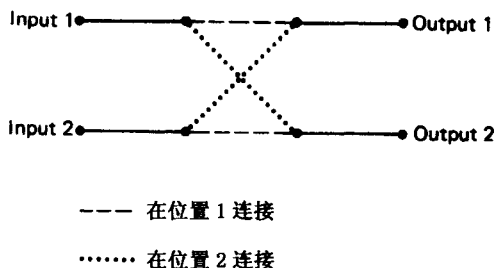


图 2-18 双位置开关

- 2.22 写一个 RASP 程序模仿以下程序计算 $\binom{n}{m}$ 。

```

procedure COMB( $n, m$ ):
  if  $m = 0$  or  $n = m$  then return 1
  else return (COMB( $n-1, m$ ) + COMB( $n-1, m-1$ ))

```

当调用时, 用堆栈存储 n 和 m 的当前值、返回值以及值地址。

- *2.23 在很多情况下, 大小为 n 的问题可以方便地分成大小为 \sqrt{n} 的 \sqrt{n} 个子问题。如下形式的递归方程结果如下

$$T\left(\frac{n^2}{2^r}\right) = nT(n) + bn^2$$

其中 r 是整数, $r \geq 1$ 。证明方程的解是 $O(n(\log n)^r \log \log n)$ 。

- 2.24 求和:

$$\begin{array}{llll}
 \text{a) } \sum_{i=1}^n i & \text{b) } \sum_{i=1}^n a^i & \text{c) } \sum_{i=1}^n ia^i & \text{d) } \sum_{i=1}^k 2^{t-i} i^2 \\
 \text{e) } \sum_{i=0}^n \binom{n}{i} & \text{f) } \sum_{i=0}^n i \binom{n}{i} & &
 \end{array}$$

- *2.25 假设 $T(1) = 1$, 解以下递推关系:

a) $T(n) = aT(n-1) + bn$

b) $T(n) = T(n/2) + bn \log n$

$$c) T(n) = aT(n-1) + bn^c$$

$$d) T(n) = aT(n/2) + bn^c$$

*2.26 通过允许递归下降到 $|S| = 1$, 修改寻找集合最大和最小元素的算法 2.3。比较次数的渐近增长率是什么?

**2.27 证明, 为了找到 n 个元素的集合中的最大和最小元素, $\lceil (3/2)n - 2 \rceil$ 次比较是必要和充分的。

*2.28 修改整数乘法算法, 分割整数为 (a)3 部分和 (b)4 部分。算法的复杂度是什么?

*2.29 设 A 是大小为 n 的正整数或负整数的数组, 其中 $A[1] < A[2] < \dots < A[n]$ 。写一个算法找到 i , 使得 $A[i] = i$, 假设这样的 i 存在。算法的执行时间复杂度为何?

证明 $O_c(\log n)$ 是最好的可能。

**2.30 如果在算法 2.4 中, n 不是 2 的幂, 将图 2-14 中的语句 $m \leftarrow (i+j-1)/2$ 替换为 $m \leftarrow \lfloor (i+j)/2 \rfloor$, 能够得到一个有效的归并排序算法。设 $T(n)$ 是用这种方法排序 n 个元素的比较次数。

a) 证明

$$T(1) = 0$$

$$T(n) = T(\lfloor n/2 \rfloor) + T(\lceil n/2 \rceil) + n - 1$$

b) 证明这个递归的解是

$$T(n) = n \lceil \log n \rceil - 2^{\lceil \log n \rceil} + 1$$

*2.31 证明递归

$$X(1) = 1$$

$$X(n) = \sum_{i=1}^{n-1} X(i)X(n-i) \quad \text{对于 } n > 1$$

的解是

$$X(n+1) = \frac{1}{n+1} \binom{2n}{n}$$

$X(n)$ 是将 n 个符号的串全部括入括号的方法数量。 $X(n)$ 的值叫做卡塔兰数。证明 $X(n) \geq 2^{n-2}$ 。

2.32 修改算法 2.5 使其输出标量乘法数量减到最小的矩阵相乘的顺序。

2.33 写出一个有效算法, 对固定的 d , 当每个 M 的维数为 1×1 , $1 \times d$, $d \times 1$ 或者 $d \times d$ 时, 能够确定计算矩阵 $M_1 \times M_2 \times \dots \times M_n$ 乘积的顺序, 使得标量乘法的数量最小。

定义 乔姆斯基范式的上下文无关语法 G 是一个 4 元组 (N, Σ, P, S) , 其中 (1) N 是一个非终结符的有限集合, (2) Σ 是一个终结符的有限集合, (3) 称 P 是一个产生的有限对集合, 具有 $A \rightarrow BC$ 或者 $A \rightarrow a$ 的形式, 其中 A, B, C 在 N 中, a 在 Σ 中, (4) S 是 N 中的特殊符号。如果 α, β, γ 是终结符和非终结符串, 而且 $A \rightarrow \beta$ 在 P 中, 则 $\alpha A \gamma \Rightarrow \alpha \beta \gamma$ 。 G 产生的语言 $L(G)$ 是终结符字符串的集合 $\{w \mid S \xRightarrow{*} w\}$, 其中 $\xRightarrow{*}$ 是 \Rightarrow 的自反和传递闭包。

2.34 写一个 $O(n^3)$ 的算法确定给定的串 $w = a_1 a_2 \dots a_n$ 是否在 $L(G)$ 中, 其中 $G = (N, \Sigma, P, S)$ 是一个乔姆斯基范式的上下文无关语法。[提示: 设 $m_{ij} = \{A \mid A \in N \text{ 和 } A \xRightarrow{} a_i a_{i+1} \dots a_j\}$ 。当且仅当 $S \in m_{1n}$ 时, $w \in L(G)$ 。可使用动态规划计算 m_{ij} 。]

*2.35 设 x 和 y 是某个字母表的符号串。考虑从 x 中删除一个符号, 在 x 中插入一个符号和替换 x 中一个符号的操作。描述一个算法, 求将 x 转换为 y 需要的最小的操作数目。

文献及注释

更多的关于数据结构及其实现的信息可在 Knuth[1968] 或 Stone[1972] 中找到。Pratt[1975]

包含用类 ALGOL 语言实现递归的描述。Berge[1958]和 Harary[1969]讨论了图论。关于树和树的遍历可参考 Knuth[1968]。Burkhard[1973]是树遍历算法的另一本参考书。

Pohl[1972]给出了算法 2.3(求最大和最小值)的优化算法。2.6 节中的 $O(n^{1.59})$ 的整数乘法算法是由 Karatsuba 和 Ofman[1962]提出的。Winograd[1973]则从更通用的角度考虑了这种加速。

动态规划的使用是由 Bellman[1957]推广开来的, 算法 2.5 是由 Godbole[1973]和 Muraoka 以及 Kuck[1973]提出的著名应用。动态规划应用于上下文无关语法的识别(习题 2.34)是 J. Cocke, Kasami[1965]和 Younger[1967]独立工作的结果。习题 2.35 的解可参阅 Wagner 和 Fischer[1974]。

更多关于解递推方程的信息, 可以参阅 Liu[1968]或者 Sloane[1973]。

第3章 排序和顺序统计

本章将考虑两个重要的相关问题,元素序列的排序和在一个序列中选择第 k 小的元素。将一个序列排序意指重新排列序列中的元素,使元素以非递增或者非递减的顺序排列。通过将序列排列为非递减序列,然后在结果序列中选择第 k 个元素,就能在序列中找到第 k 小的元素。不过,我们会看到,还有比这种方法更快地选择第 k 小元素的方法。

排序是一个对于实践很重要同时在理论上又很有趣的问题。因为相当一部分商业数据的处理都涉及对大量数据排序,所以,有效的排序算法对于经济性而言相当重要。即使在算法设计中,对元素序列的排序处理也是很多算法的基本部分。

本章将考虑两类排序算法。第一类算法利用元素的结构进行排序。例如,如果排序元素是在固定范围 $0 \sim m-1$ 中的整数,那么能够在 $O(n+m)$ 的时间内对 n 个元素的序列排序。如果要排序的元素是固定字母表的字符串,那么字符串序列的排序时间与字符串长度的和是线性成比例的。

第二类算法假设待排序的元素没有结构。基本操作是一对元素间的比较。由于算法具有这种特性,可以看出对有 n 个元素的序列进行排序至少需要 $n \log n$ 次比较。这里将给出两个 $O_c(n \log n)$ 排序算法:一种是堆排序,最坏情况下复杂度是 $O_c(n \log n)$;另一种是快速排序,期望情况下时间复杂度是 $O_c(n \log n)$ 。

3.1 排序问题

定义 集合 S 的偏序是一个关系 R ,对 S 中的元素 a 、 b 和 c 均有:

- 1) aRa 为真(R 是自反的)
- 2) aRb 和 bRc 蕴涵 aRc (R 是传递的)
- 3) aRb 和 bRa 蕴涵 $a=b$ (R 是反对称的)

整数的关系 \leq 和关系 \subseteq (集合包含)是偏序的两个例子。

集合 S 的线性序或者全序是 S 上的一个偏序关系,使得对每对元素 a, b ,有 aRb 或者 bRa 。整数的关系 \leq 是一个线性序;[⊙]集合的关系 \subseteq 不是线性序。

可以用如下方式表述排序问题。给定一个集合的 n 个元素的序列 a_1, a_2, \dots, a_n ,该集合上的线性序经常表示为 \leq 。我们可以找出这 n 个元素的一个排列 $\pi: a_{\pi(1)}, a_{\pi(2)}, \dots, a_{\pi(n)}$,它将给定序列映射为一个非递减序列,使 $a_{\pi(i)} \leq a_{\pi(i+1)}$, $1 \leq i < n$ 。通常需要对序列自身进行排序而不是对排列 π 排序。

排序方法分为内部排序(数据已经在随机访问存储器中)和外部排序(数据主要在随机访问存储器之外)。外部排序作为计算数据处理应用的一部分,通常涉及的元素数量比在随机访问存储器中进行一次排序的元素多很多。因此,在辅助存储设备(如磁盘或者磁带)上进行数据的外部排序具有重大的商业意义。

内部排序在算法设计和商业应用中都很重要。在排序作为其他算法的一部分的情况下,为

⊙ 对于线性序 \leq ,正如通常期望的那样,可用 $a < b$ 表示 $a \leq b$ 且 $a \neq b$ 。而且, $b > a$ 与 $a < b$ 同义, $b \geq a$ 也与 $a \leq b$ 同义。

了适合随机访问存储器,要排序的数据项的数量通常比较小。然而,这里假设要排序的数据项的数量足够。如果只是对少量数据项排序,使用简单的策略,如复杂度为 $O(n^2)$ 的“冒泡排序”(见习题 3.5)更合适。

排序算法有很多种。本书并不想列举所有重要的排序算法,而只介绍在算法设计中有用的方法。首先考虑要排序的元素是整数或者(几乎等同于)有限字母表的字符串的情况。可以看到,此类排序能够在线性时间内完成。接下来,考虑没有利用整数和字符串特殊性质的排序问题,此时,只能依赖于被排序元素的比较结果对程序进行分支。在这种情况下,对有 n 个元素的序列排序, $O(n \log n)$ 次比较是必要的也是充分的。

3.2 基数排序

现在开始研究整数排序,令 a_1, a_2, \dots, a_n 是 $0 \sim m-1$ 之间的整数组成的序列。如果 m 不是太大,这个序列可以很容易地用以下方式排序。

- 1) 初始化 m 个空队列,每一个队列对应 $0 \sim m-1$ 中的一个整数。每个队列称为桶。
- 2) 从左到右扫描序列 a_1, a_2, \dots, a_n , 将元素 a_i 放到第 a_i 个队列中。
- 3) 连接这些队列(队列 $i+1$ 的内容加到队列 i 的末端), 得到排序结果。

由于能够在常数时间内将一个元素插入第 i 个队列,所以 n 个元素能够在 $O(n)$ 时间内插入队列。连接 m 个队列需要的时间为 $O(m)$ 。如果 m 是 $O(n)$, 那么算法用 $O(n)$ 的时间排序 n 个整数。这种算法一般称为桶排序。

桶排序能够扩展为对整数数组(例如表)序列进行排序,此时按字典序排序。设 \leq 为集合 S 的线性序,如果 $(s_1, s_2, \dots, s_p) \leq (t_1, t_2, \dots, t_q)$, 确切地说,当以下两个条件之一满足时:

- 1) 存在整数 j , 使得 $s_j < t_j$, 而且对于所有的 $i < j$, $s_i = t_i$ 。
- 2) 对于 $1 \leq i \leq p$, $p \leq q$ 而且 $s_i = t_i$ 。

被扩展为元组(其元素来自 S)的关系 \leq 是字典序。例如,如果将字母组成的字符串(基于自然字母顺序)当作元组,那么字典中的词就是以字典序排列。

首先将桶排序推广到由 k 元组序列, k 元组的元素是 $0 \sim m-1$ 间的整数。通过 k 次对序列进行桶排序操作,可实现排序。在第一次操作时, k 元组根据它们的第 k 个元素排序。第二次操作的结果序列根据第 $(k-1)$ 个元素排序。根据第二次操作的结果,第三次操作序列对第 $(k-2)$ 个元素排序,依此类推。第 k 次(最后一次)操作序列根据第一个元素对 $(k-1)$ 次的结果排序。[⊙] 序列现在是字典序。算法的准确描述如下。

算法 3.1 字典排序。

输入: 一个序列 A_1, A_2, \dots, A_n , 其中 A_i 是一个 k 元组

$$(a_{i1}, a_{i2}, \dots, a_{ik})$$

a_{ij} 是一个 $0 \sim m-1$ 之间的整数。(实现 k 元组序列的便捷数据结构是一个 $n \times k$ 的数组。)

输出: 一个序列 B_1, B_2, \dots, B_n , 它是 A_1, A_2, \dots, A_n 的排列, 满足 $B_i \leq B_{i+1}$, 其中 $1 \leq i < n$ 。

方法: 将 k 元组 A_i 传递给某个桶, 只需移动指向 A_i 的指针。这样, A_i 能够在固定的时间而不是在 k 时间内添加到一个桶中。用称为 QUEUE 的队列保存“当前”的元素序列。用数组 Q 来存储 m 个桶, 其中桶 $Q[i]$ 存储当前操作元素中包含着整数 i 的 k 元组。算法如图 3-1 所示。 □

⊙ 在很多情况下, 仅根据第一个元素进行字符串的桶排序就足够了。如果放入桶的元素数量很小, 那么可以直接如冒泡排序, 对每个桶中的字符串进行排序。

```

begin
  place  $A_1, A_2, \dots, A_n$  in QUEUE;
  for  $j \leftarrow k$  step  $-1$  until  $1$  do
    begin
      for  $l \leftarrow 0$  until  $m-1$  do make  $Q[l]$  empty;
      while QUEUE not empty do
        begin
          let  $A_i$  be the first element in QUEUE;
          move  $A_i$  from QUEUE to bucket  $Q[a_{ij}]$ 
        end;
      for  $l \leftarrow 0$  until  $m-1$  do
        concatenate contents of  $Q[l]$  to the end of QUEUE
      end
    end
  end
end

```

图 3-1 字典排序算法

定理 3.1 算法 3.1 在 $O((m+n)k)$ 时间内, 按字典序对长度为 n 的 k 元组序列进行排序, 其中, k 元组的每个元素都是 $0 \sim m-1$ 的整数。

证明: 通过对外部循环的执行次数进行归纳, 以证明算法 3.1 是正确的。归纳假设是在 r 次执行后, QUEUE 中的 k 元组已按照右边第 r 个元素的字典序排好。一旦观察到第 $(r+1)$ 次排序是对 k 元组右边的第 $(r+1)$ 个元素进行操作的事实, 如果两个 k 元组放在同一个桶中且按照右边第 r 个元素的字典序排列。第一个 k 元组位于第二个 k 元组之前, 则可以很容易得出结果。

算法 3.1 的一趟外部循环需要的时间为 $O(m+n)$ 。循环重复 k 次的时间复杂度为 $O((m+n)k)$ 。□

算法 3.1 有各种应用。它在很长一段时间内用于穿孔卡片排序机。它也可以用于在 $O(kn)$ 的时间内排序 $O(n)$ 个 $0 \sim n^k-1$ 之间的整数, 因为可以将这种整数看作以 $0 \sim n-1$ 之间的数组成的 k 元组。(即以 n 为基数的整数表示。)

桶排序最终推广到应用于大小不同的元组(称为串)。如果最长的串长度为 k , 则可以给每个串填充一个特殊的符号, 使得所有的串长度为 k , 然后应用算法 3.1。然而, 如果只有几个长串, 则这种方法会由于以下两个原因而造成效率低下。首先, 每次操作都要检测每个串, 其次, 即使几乎所有的桶都是空的, 还是要检测每个桶 $Q[i]$ 。下面将给出一种时间为 $O(m+l_{\text{total}})$, 对 n 个不同长度的串序列排序的算法, 其中 l_i 是第 i 个串的长度, 而且 $l_{\text{total}} = \sum_{i=1}^n l_i$ 。串的元素在 $0 \sim m-1$ 的范围内。该算法在 m 和 l_{total} 都是 $O(n)$ 时, 是很有用的。

该算法的本质是首先将串按长度递减的顺序排序。设 l_{max} 为最长的串的长度。接下来, 像算法 3.1 那样, 进行 l_{max} 次桶排序操作。然而, 第一次排序操作(按最右边的元素排序)只对长度为 l_{max} 的串进行。第二次排序(根据第 $(l_{\text{max}}-1)$ 个元素)时串的长度至少为 $l_{\text{max}}-1$, 依此类推。

例如, 假设 bab , abc 和 a 是三个待排序的串(虽然假设元组的元素是整数, 但为了表示方便, 我们使用字母。这不会带来困难, 因为可以用 0 、 1 和 2 来代替 a 、 b 和 c)。这里 $l_{\text{max}}=3$, 所以, 第一次排序仅对前两个串按照它们的第三个元素进行排序。在排序时, 将 bab 放入 b 桶, 将 abc 放入 c 桶。 a 桶为空。第二次操作则根据前两个串的第二个元素对它们进行排序。现在, a 桶和 b 桶将放入串, c 桶仍然为空。在第三次和最后一次排序中, 对所有三个串根据它们的第一个元素进行排序。这次 a 桶和 b 桶将放入串, c 桶为空。

可以看到, 通常情况下, 某一次排序时很多桶可能为空。这样, 在某一次排序时判定哪些桶是非空的处理步骤是很有益的。由于每一趟的非空桶列表由递增桶编号确定, 这允许算法在与非空桶数成比例的时间内对非空桶进行连接。

算法 3.2 对不同长度的串进行字典排序。

输入：串(元组)序列 A_1, A_2, \dots, A_n ，其元素是 $0 \sim m-1$ 之间的整数。设 l_i 是 $A_i = (a_{i1}, a_{i2}, \dots, a_{il_i})$ 的长度，设 l_{\max} 是 l_i 中的最大值。

输出： A_i 的排列 B_1, B_2, \dots, B_n ，满足 $B_1 \leq B_2 \leq \dots \leq B_n$ 。

方法：

1) 对每个 $l (1 \leq l \leq l_{\max})$ ，为符号造一个表，这些符号是出现在一个或多个串的第 l 个元素中的。首先对 A_i 的每个元素 a_{il} ，其中 $1 \leq i \leq n, 1 \leq l \leq l_i$ ，创建一个 (l, a_{il}) 对。这样的对表示某串的第 l 个元素包含整数 a_{il} ；然后使用算法 3.1 的显式推广进行字典排序。^① 接下来，通过从左到右扫描排序表，很容易创建 l_{\max} 个排序表 $\text{NONEMPTY}[l]$ ，其中 $1 \leq l \leq l_{\max}$ ，这样， $\text{NONEMPTY}[l]$ 确实包含出现在串的第 l 个元素的符号。也就是说， $\text{NONEMPTY}[l]$ 有序地包含了所有整数 j ，使得对某些 $i, a_{il} = j$ 。

2) 确定每个串的长度。创建表 $\text{LENGTH}[l]$ ， $1 \leq l \leq l_{\max}$ ，其中 $\text{LENGTH}[l]$ 包括所有长度为 l 的串。(虽然我们说的是移动一个串，但实际上移动的仅仅是指向串的指针。这样，每个串都能够用固定数量的步骤加入到 $\text{LENGTH}[l]$ 中。)

3) 现在，像算法 3.1 那样，从位置 l_{\max} 处的元素开始对串排序。然而，在第 i 次排序后， QUEUE 只包含长度为 $l_{\max} - i + 1$ 或更大的串，这些串的 $l_{\max} - i + 1$ 到 l_{\max} 元素已经按字典序排好了。在步骤 1 中计算的表 NONEMPTY 用来确定在每次桶排序操作中占据哪些桶。这个信息有助于加快桶的连接速度。图 3-2 给出了算法的简化 ALGOL 语言实现。□

```

begin
1.   make QUEUE empty;
2.   for j ← 0 until m-1 do make Q[j] empty;
3.   for l ← lmax step -1 until 1 do
      begin
4.         concatenate LENGTH[l] to the beginning of
           QUEUE;①
5.         while QUEUE not empty do
              begin
6.                 let Al be the first string on QUEUE;
7.                 move Al from QUEUE to bucket Q[al]
              end;
8.         for each j on NONEMPTY[l] do
              begin
9.                 concatenate Q[j] to the end of QUEUE;
10.                make Q[j] empty
              end
            end
      end
end

```

① 从技术上说，仅连接到队列的末尾，但连接到队列头并从概念上说并不会带来困难。针对本例，最有效的方法是首先从第 6 行的 $\text{LENGTH}[l]$ 中选择 A_l 的值，然后从 QUEUE 中选择它们，而不用连接 $\text{LENGTH}[l]$ 和 QUEUE 。

图 3-2 不同长度串的字典排序

例 3.1 下面用算法 3.2 排序串 a , bab 和 abc 。图 3-3 是一种可能的表示这些串的数据结构。 STRING 是一个数组， $\text{STRING}[i]$ 是指向第 i 个串的表示的指针，串的长度和元素存储在数组 DATA 中。 $\text{STRING}[i]$ 指向的 DATA 中的单元给出第 i 个串的符号数 j 。 DATA 中接下来的 j 个单元包含这些符号。

① 在算法 3.1 中，假设元素是从同一个字母表中选择的。这里第二个元素的取值范围是 $0 \sim m-1$ ，而第一个元素的取值范围是 $1 \sim l_{\max}$ 。

和算法 3.2 所用的串表实际上是指针表, 就像数组 STRING 那样。为了表示方便, 本例后面的部分将在表中写入实际的串, 而不是指向串的指针。记住, 存储在队列中的是指针而不是串本身。

算法 3.2 的第 1 部分针对第一个串创建了对 (1, *a*); 第二个串创建对 (1, *b*), (2, *a*), (3, *b*); 第三个串创建对 (1, *a*), (2, *b*), (3, *c*)。这些对的有序表是:

(1, *a*), (1, *a*), (1, *b*), (2, *a*), (2, *b*), (3, *b*), (3, *c*)

从左到右扫描这个排序表, 可推出

NONEMPTY[1] = *a*, *b*

NONEMPTY[2] = *a*, *b*

NONEMPTY[3] = *b*, *c*

算法 3.2 的第 2 部分计算 $l_1 = 1$, $l_2 = 3$ 和 $l_3 = 3$ 。于是, LENGTH[1] = *a*, LENGTH[2] 为空, LENGTH[3] = *bab*, *abc*。因此, 在算法第 3 部分一开始就设 QUEUE = *bab*, *abc* 并按照它们的第三个元素排序这些串。

NONEMPTY[3] = *b*, *c* 使我们确信当在图 3-2 的第 8~10 行生成排序表时, $Q[a]$ 不需要连接在 QUEUE 的尾部。这样, 在图 3-2 的第 3~10 行的第一次循环执行完毕之后, QUEUE = *bab*, *abc*。

在第二次操作时, QUEUE 不变, 因为 LENGTH[2] 为空, 而且按照第二个元素排序没有改变顺序。第三次操作时, 在第 4 行设 QUEUE 为 *a*, *bab*, *abc*。根据第一个元素排序使得 QUEUE = *a*, *abc*, *bab*, 这就是正确的顺序。注意, 第三次操作时, $Q[c]$ 保持为空, 因为 *c* 不在 NONEMPTY[1] 中, 因此不将 $Q[c]$ 连接到 QUEUE 尾部。□

定理 3.2 算法 3.2 在时间 $O(l_{\text{total}} + m)$ 内对输入进行排序, 其中

$$l_{\text{total}} = \sum_{i=1}^n l_i$$

证明: 通过对图 3-2 的外部循环操作次数的简单归纳, 证明在第 i 次操作后, QUEUE 包含长度大于等于 $l_{\text{max}} - i + 1$ 的串, 根据 $l_{\text{max}} - i + 1$ 到 l_{max} 的元素排序。这样, 算法按字典序对它的输入进行排序。

从用时结果看, 算法的第 1 部分用 $O(l_{\text{total}})$ 时间创建对, 用 $O(m + l_{\text{total}})$ 时间排序它们。类似地, 第 2 部分需要的时间不超过 $O(l_{\text{total}})$ 。

现在必须将注意力转向算法的第 3 部分和图 3-2 所示的程序。设 n_i 是有第 i 个元素的串的数量, m_i 是出现在串的第 i 个元素中的不同符号的数量 (即 m_i 是 NONEMPTY[i] 的长度)。

考虑图 3-2 第 3 行中 l 的一个固定值。5~7 行的循环需要的时间为 $O(n_l)$, 8~10 行的循环需要的时间为 $O(m_l)$ 。第 4 步需要的时间为一个常数, 因此, 通过 3~10 行的一次操作需要的时间为 $O(m_l + n_l)$ 。因此, 整个循环需要的时间为:

$$O\left(\sum_{l=1}^{l_{\text{max}}} (m_l + n_l)\right)$$

由于

$$\sum_{l=1}^{l_{\text{max}}} m_l \leq l_{\text{total}} \text{ 和 } \sum_{l=1}^{l_{\text{max}}} n_l = l_{\text{total}}$$

可以看到 3~10 行需要的时间为 $O(l_{\text{total}})$ 。由于第 1 行需要常数时间, 第 2 行需要的时间为 $O(m)$, 就得到了所期望的结果。□

下面给出一个例子, 说明如何在算法设计中实现对串进行升序排序。

例 3.2 如果能够通过改变顶点的儿子的顺序使一棵树映射为另一棵树, 则称这两棵树是同

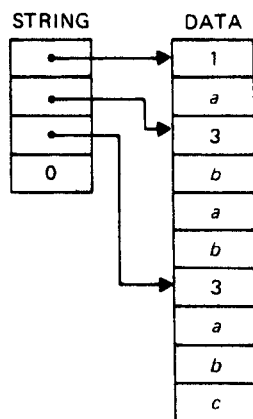


图 3-3 串的数据结构

构的。考虑判断两棵树 T_1 、 T_2 是否同构的问题。以下算法所需时间与结点数成线性比例。算法将整数赋给两棵树的结点, 从第 0 级的结点开始, 向上操作至根, 在这种方式下, 当且仅当给它们的根赋同样的整数时, 两棵树是同构的。算法的执行过程如下:

1) 给 T_1 、 T_2 的所有叶子赋值为整数 0。

2) 归纳假设, 已经对 T_1 、 T_2 的 $i-1$ 级的所有结点赋整数值。假设 L_1 是 T_1 的第 $i-1$ 级的结点表, 该表根据所赋的整数值以非递减顺序排列, L_2 是 T_2 的对应表。[○]

3) 通过从左到右扫描表 L_1 并执行以下操作给 T_1 的第 i 级非叶结点赋予一个整数元组: 取 L_1 的每个结点 v 所赋值的整数作为与 v 的父结点相关联的元组的下一个元素。在这一步的最后, T_1 中第 i 级的每个非叶结点 w 将有元组 (i_1, i_2, \dots, i_k) 与其相关联, 其中 i_1, i_2, \dots, i_k 是与 w 的儿子结点相关联的按非递减顺序排列的整数。设 S_1 是为 T_1 的第 i 级结点创建的元组序列。

4) 对 T_2 重复步骤 3, 设 S_2 是为 T_2 的第 i 级结点创建的元组序列。

5) 用算法 3.2 排序 S_1 和 S_2 。设 S'_1 和 S'_2 分别为已排序的元组序列。

6) 如果 S'_1 和 S'_2 不相同, 则停止执行且判定树不是同构的。否则, 将整数 1 赋值给用 S'_1 的第一个不同元组表示的 T_1 第 i 级的结点, 将整数 2 赋给第二个不同元组表示的结点, 依此类推。将这些整数赋给 T_1 第 i 级的结点后, 创建一个赋值结点的表 L_i 。将 T_1 第 i 级的所有叶子结点添加到表 L_i 的前面。设 L_2 为与 T_2 的结点对应的表。返回步骤 3, 用这两个表给 $i+1$ 级的结点元组赋值。

7) 如果给 T_1 和 T_2 的根赋予了相同的整数, 则它们是同构的。□

图 3-4 对两棵同构树的结点进行整数和元组赋值的结果。

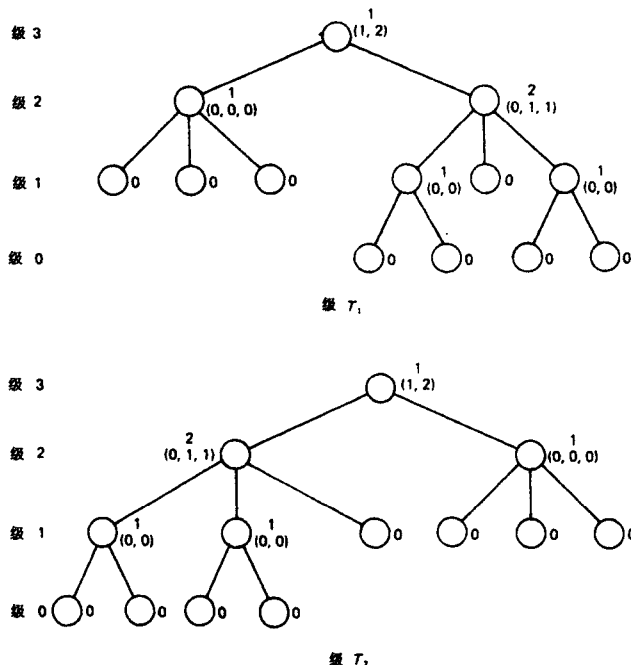


图 3-4 同构树算法的赋值

定理 3.3 可以在 $O(n)$ 时间内确定两棵有 n 个结点的树是否同构。

证明: 定理由对例 3.2 的算法进行形式化而得到。算法的正确性证明略。通过观察将整数赋

○ 必须确保能够通过前序遍历在 $O(n)$ 步内进行级别数赋值。

给第 i 级的结点(而不是叶子)的工作是与 $i-1$ 级的结点数成比例的,从而分析出运行时间。将所有级的时间相加得到时间 $O(n)$ 。将整数赋给叶子的时间也与 n 成比例,这样,算法所用时间为 $O(n)$ 。□

结点都附有标记的树是标记树。假设顶点标记是 $1 \sim n$ 之间的整数。如果将每个结点的标记作为根据上述算法赋给结点的元组的第一个元素,就能够在线性时间内确定两棵有 n 个结点的标记树是否同构。这样,就得到以下推论。

推论 判断两棵标记范围为 $1 \sim n$ 的有 n 个结点的标记树是否同构需要 $O(n)$ 时间。

3.3 比较排序

这一节考虑对从线性序集合 S 中抽出 n 个元素进行排序,其中 S 的元素的结构未知。用于获得关于序列的信息的唯一操作是比较两个元素。本节先说明任何基于比较的对长度为 n 的序列进行排序的算法至少需要进行 $O(n \log n)$ 次比较。

假设有 n 个不同的待排序元素 a_1, a_2, \dots, a_n 。比较排序的算法可以用 1.5 节描述的决策树表示。图 1-18 给出了一棵为 a, b, c 排序的决策树。在决策树的某个结点 v 比较元素 a 与元素 b , 如果 $a < b$, 那么分支转向 v 的左儿子, 如果 $a \geq b$, 那么分支转向 v 的右儿子。

通常, 基于比较的排序算法限制一次只能比较两个输入元素。但实际上, 应用于任意线性序集合的算法并不以任何形式与输入数据绑定在一起。因为在一般情况下, 对数据操作并“没有意义”。在任何情况下, 都能确切证明对 n 个元素进行排序的决策树的高度。

引理 3.1 高度为 h 的二叉树最多有 2^h 个叶子。

证明: 对 h 进行基本归纳。只需要注意, 高度为 h 的二叉树由一个根和至多两棵子树组成, 每棵子树的高度最多为 $h-1$ 。□

定理 3.4 任何排序 n 个不同元素的决策树的高度至少为 $\log n!$ 。

证明: 因为排列 n 个元素可能的结果是输入的 $n!$ 个排列中的任何一个, 所以决策树肯定至少有 $n!$ 个叶子。由引理 3.1 可得, 高度必须至少为 $\log n!$ 。□

推论 任何比较排序算法至少需要 $cn \log n$ 次比较来排序 n 个元素, 其中 $c > 0$, n 足够大。

证明: 对 $n > 1$

$$n! \geq n(n-1)(n-2) \cdots \left(\lceil \frac{n}{2} \rceil\right) \geq \left(\frac{n}{2}\right)^{n/2}$$

所以, 对 $n \geq 4$, $\log n! \geq (n/2) \log(n/2) \geq (n/4) \log n$ 。□

用斯特林近似对 $n!$ 进行更精确的估计, 值为 $(n/e)^n$, 所以 $n(\log n - \log e) = n \log n - 1.44n$ 接近排序 n 个元素需要的比较次数的下界, 是一个令人满意的近似值。

3.4 堆排序—— $O(n \log n)$ 的比较排序算法

要对长度为 n 的序列排列, 每个基于比较的排序算法本质上需要 $n \log n$ 次比较, 很自然存在这样一个疑问, 是否存在只用 $O(n \log n)$ 次比较就能对长度为 n 的序列进行排序的算法? 实际上, 前面已经讨论过这样的算法, 即 2.7 节的归并排序。另一个这种算法是堆排序, 除了是有用的排序算法外, 堆排序使用的有趣的数据结构还有许多其他用途。

对堆排序最好的理解是考虑图 3-5 中的二叉树, 它的每片树叶的深度为 d 或者 $d-1$ 。可用待排序的序列的元素标记树的结点。堆排序重新安排树的元素, 直到与顶点相关的元素大于等于与该结点的儿子相关的元素为止。这样的标记树称为堆。

例 3.3 图 3-5 给出了一个堆。注意, 从每片叶子到根的路径上的元素序列是线性有序的, 子树中的最大元素总是那棵子树的根。□

堆排序的下一步是从堆中移走最大的元素, 该元素通常位于根部。将某片叶子的标记移到根, 并将该叶子删除。这样, 树重新成为一个堆, 并重复这个过程。从堆移除的元素序列就是已排好序的元素序列(以降序排列)。

对于堆来说, 一个方便的数据结构是数组 A , 其中 $A[1]$ 存储根元素, $A[2i]$ 和 $A[2i+1]$ 存储元素 $A[i]$ 所对应的顶点的左右孩子处的元素(如果存在的话)。例如, 图 3-5 中的堆可用如下数组表示:

16	11	9	10	5	6	8	1	2	4
----	----	---	----	---	---	---	---	---	---

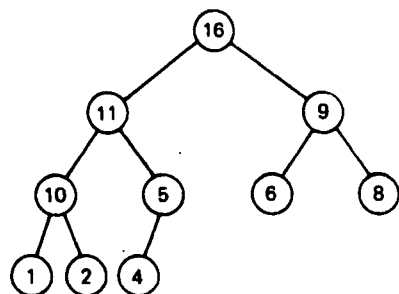


图 3-5 堆

可以看出, 深度最小的顶点在数组中首先出现, 深度相等的顶点以从左到右的顺序出现。

并不是每个堆都可以用此方式表示。如图 3-5 所示, 按照树的表示法, 最低层的叶子必须尽可能地靠左放置。

如果用数组表示堆, 那么堆排序算法的几个操作是很容易执行的。例如, 在算法中, 必须将根元素移除, 并将这个元素存储到其他地方, 然后使剩下的树重新组成堆, 同时移去未标记的叶子。我们可以移除堆的最大元素, 并通过交换 $A[1]$ 和 $A[n]$ 来存储它。然后, 认为数组存储单元 n 不再是堆的一部分。存储单元 n 作为从树删除的叶子。为了使树的存储单元 $1, 2, \dots, n-1$ 重新形成堆, 在加入新的元素 $A[1]$ 后, 按照需要从树根沿路径向下移动。接下来重复这个过程, 交换 $A[1]$ 和 $A[n-1]$, 并认为树占据的存储单元为 $1, 2, \dots, n-2$, 依此类推。

例 3.4 以图 3-5 的堆为例, 考虑当交换表示堆的数组的第一个和最后一个元素的时候发生了什么。结果数组对应于如图 3-6a 所示的标记树。

4	11	9	10	5	6	8	1	2	16
---	----	---	----	---	---	---	---	---	----

进一步考虑排除元素 16。为了将结果树转变为堆, 将元素 4 与 11 交换, 即它的两个孩子 9 和 11 中较大的那个孩子 11。

在新位置, 4 的孩子是 10 和 5。因为它们都比 4 大, 因此交换 4 和它较大的孩子 10。这样, 在 4 的新位置, 它的孩子是 1 和 2。因为 4 比它的每个孩子都大, 所以不需要进一步交换。结果堆如图 3-6b 所示。注意, 虽然元素 16 已经从堆中移除, 它仍然存储于数组 A 的尾部。□

现在开始正式描述堆排序算法。设 a_1, a_2, \dots, a_n 是待排序的元素序列。假设元素在数组 A 中初始顺序为 $A[i] = a_i, 1 \leq i \leq n$ 。第一步是创建堆, 也就是说, 重新排列数组 A 中的元素以满足堆的性质: 对 $1 \leq i \leq n/2, A[i] \geq A[2i]$; 对于 $1 \leq i < n/2$, 有 $A[i] \geq A[2i+1]$ 。过程是由叶子开始逐渐构建越来越大的堆。每棵包含一片叶子的子树形成一个堆。如果根部元素小于其儿子的元素, 那么高度为 h 的子树通过将根部的元素和儿子元素中较大的那个进行交换来形成堆。但这样做可能会破坏高度为 $h-1$ 的堆, 该高度为 $h-1$ 的堆需要重新构建。详细算法如下所示。

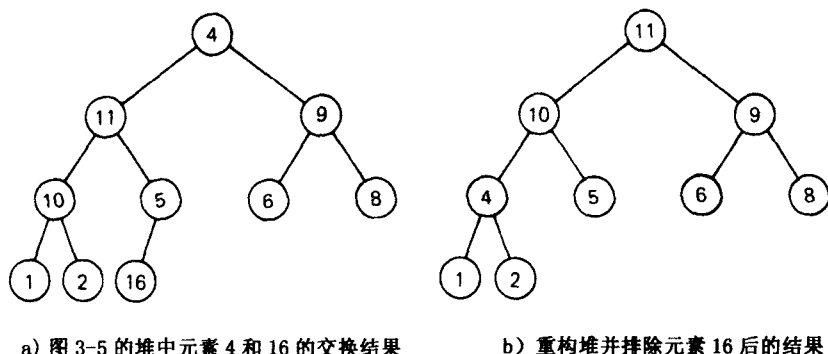


图 3-6

算法 3.3 构建堆。

输入：元素数组 $A[i]$, $1 \leq i \leq n$ 。

输出：已构建为堆的 A 的元素，满足对 $1 < i \leq n$, $A[i] \leq A[\lfloor i/2 \rfloor]$ 。

方法：算法的核心是递归过程 HEAPIFY。参数 i 和 j 给出数组 A 具有堆性质的存储单元的范围，其中 i 是即将要创建的堆的根。

procedure HEAPIFY(i, j):

1. **if** i is not a leaf and if a son of i contains a larger element than i
 do then
 begin
2. let k be a son of i with the largest element;
3. interchange $A[i]$ and $A[k]$;
4. HEAPIFY(k, j)
- end**

参数 j 用于确定 i 是否为叶子和 i 是否有一个或者两个孩子。如果 $i > j/2$ ，那么 i 是叶子，HEAPIFY(i, j) 不需要执行任何操作，因为 $A[i]$ 本身就是堆。

使所有 A 具备堆的性质的算法是很简单的：

procedure BUILDHEAP;

for $i \leftarrow n \ominus \text{step} - 1$ **until** 1 **do** HEAPIFY(i, n)

□

接下来，我们说明算法 3.3 在线性时间内以 A 构建一个堆。

引理 3.2 如果顶点 $i+1, i+2, \dots, n$ 是堆的根，那么在调用 HEAPIFY(i, n) 后，所有的 $i, i+1, \dots, n$ 都将是堆的根。

证明：通过对 i 向后归纳进行证明。归纳基础 $i=n$ 是平凡的，因为顶点 n 必须是叶子，而且第 1 行的测试确定 HEAPIFY(n, n) 什么都没有做。

对于归纳步骤，要注意，如果顶点 i 是叶子或者其孩子没有比它大的元素，那么根据以上分析，就不需要证明了。然而，如果顶点 i 有一个孩子（即 $2i=n$ ），而且 $A[i] < A[2i]$ ，那么第 3 行的 HEAPIFY 将交换 $A[i]$ 和 $A[2i]$ 。在第 4 行，调用 HEAPIFY($2i, n$)，所以归纳假设意味着将以顶点 $2i$ 为根的树重构为堆。顶点 $i+1, i+2, \dots, 2i-1$ 从来没有停止作为堆的根。因为在数组 A 的新排列中有 $A[i] > A[2i]$ ，所以 i 为根的树也是堆。

同样，如果顶点 i 有两个孩子（即 $2i+1 \leq n$ ），而且 $A[2i]$ 和 $A[2i+1]$ 中较大的一个比 $A[i]$ 大，如上所示，可以知道在调用 HEAPIFY(i, n) 后，所有的 $i, i+1, \dots, n$ 会是堆的根。 □

⊖ 实际上，以 $\lfloor n/2 \rfloor$ 开始。

定理 3.5 算法 3.3 使 A 在线性时间内成为一个堆。

证明：使用引理 3.2，对 $1 \leq i \leq n$ ，可以通过对 i 进行简单的向后归纳说明顶点 i 成为堆的根。

设 $T(h)$ 是在高度为 h 的顶点上执行 HEAPIFY 需要的时间，那么对某个常数 c ， $T(h) \leq T(h-1) + c$ ，说明 $T(h)$ 是 $O(h)$ 。

算法 3.3 对每个顶点调用一次 HEAPIFY，排除了递归调用。这样，BUILDHEAP 花费的时间与所有顶点的高度总和同阶。但是，最多有 $\lceil n/2^{i-1} \rceil$ 个顶点高度为 i 。因此，BUILDHEAP 花费的总时间与

$$\sum_{i=1}^h i \frac{n}{2^i}$$

同阶，即 $O(n)$ 。 □

现在可以完成堆排序的说明。一旦 A 的元素排序为堆，每次从根移除一个元素。该过程通过交换 $A[1]$ 和 $A[n]$ ，并重新排列 $A[1]$ ， $A[2]$ ， \dots ， $A[n-1]$ 为堆来实现。接下来，交换 $A[1]$ 和 $A[n-1]$ ，并重新排列 $A[1]$ ， $A[2]$ ， \dots ， $A[n-2]$ 为堆，依此类推，直到堆只包含一个元素为止。到那个时候， $A[1]$ ， $A[2]$ ， \dots ， $A[n]$ 就是已排好序的元素序列了。

算法 3.4 堆排序。

输入：数组元素 $A[i]$ ， $1 \leq i \leq n$ 。

输出：按非递减顺序排列的 A 的元素。

方法：使用过程 BUILDHEAP，即算法 3.3。算法如下：

```
begin
  BUILDHEAP;
  for  $i \leftarrow n$  step  $-1$  until 2 do
    begin
      interchange  $A[1]$  and  $A[i]$ ;
      HEAPIFY(1,  $i-1$ )
    end
  end
```

□

定理 3.6 算法 3.4 在时间 $O(n \log n)$ 内排序 n 个元素。

证明：通过对主循环的执行次数 m 进行归纳，证明算法的正确性。归纳假设是在 m 次迭代后， $A[n-m+1]$ ， \dots ， $A[n]$ 包含 m 个已排序的最大元素（最小元素在前），并且 $A[1]$ ， \dots ， $A[n-m]$ 构成堆。具体的证明留作练习。执行 HEAPIFY(1, i) 的时间是 $O(\log i)$ 。因此算法 3.4 的复杂度是 $O(n \log n)$ 。 □

推论 堆排序时间复杂度为 $O_c(n \log n)$ 。

3.5 快速排序——期望时间为 $O(n \log n)$ 的排序算法

到目前，仅考虑了最坏情况下的排序算法。对许多应用，算法时间复杂度更实际的测量是期望运行时间。在决策树模型下考察排序，可以看到没有一个基于比较的排序算法的期望时间复杂度能够明显低于 $n \log n$ 。然而，对某个常数 c ，确实可以找到最坏运行时间为 cn^2 ，但是期望运行时间是已知排序算法中最好的算法。本节讨论的快速排序算法就是这种算法的一个例子。

在讨论算法的期望运行时间之前，必须统一输入的概率分布。对排序来说，一个自然的、

也是我们要做的假设是待排序序列的每个排列在输入中出现可能性是相等的。在这个假设下,我们能够确定排序 n 个元素所需的期望比较次数。

一般的方法是将给定输入达到叶子 v 的概率与决策树的每片叶子 v 联系起来。如果知道输入的概率分布,就可以确定与叶子相连的概率。从而可以估算对应特定排序算法的决策树的所有叶子的比较次数的期望值的总和 $\sum p_i d_i$, 其中 p_i 是到达第 i 片叶子的概率, d_i 是它的深度。称这个数值为决策树的期望深度。下面是定理 3.4 的推广。

定理 3.7 在所有 n 个元素的排列作为输入出现的可能性相等的假设下, 排序 n 个元素的决策树的期望深度至少为 $\log n!$ 。

证明: 设 $D(T)$ 是二叉树 T 的叶子深度的和。设 $D(m)$ 是有 m 片叶子的二叉树 T 的 $D(T)$ 的最小值。通过对 m 做归纳, 可证明 $D(m) \geq m \log m$ 。

归纳基础 $m=1$ 是平凡的。现在, 假设对所有值小于 k 的 m 归纳假设设为真, 考虑有 k 片叶子的决策树 T 。对某个 i , $1 \leq i < k$, T 的根包括一棵有 i 片叶子的左子树 T_i 和一棵有 $k-i$ 片叶子的右子树 T_{k-i} 。显然,

$$D(T) = i + D(T_i) + (k-i) + D(T_{k-i})$$

因此, 最小的和可通过下式给出

$$D(k) = \min_{1 \leq i \leq k} [k + D(i) + D(k-i)] \quad (3-1)$$

使用归纳假设, 从式(3-1)可以得到

$$D(k) \geq k + \min_{1 \leq i \leq k} [i \log i + (k-i) \log(k-i)] \quad (3-2)$$

容易看出最小值出现在 $i=k/2$ 处。这样,

$$D(k) \geq k + k \log \frac{k}{2} = k \log k$$

由此得到, 对所有 $m \geq 1$, $D(m) \geq m \log m$ 。

现在说明排序 n 个随机元素的决策树 T 有至少 $n!$ 片叶子。确切地说, $n!$ 片叶子的概率都为 $1/n!$, 剩下的叶子的概率为零。对于祖先皆是概率为 0 的叶子顶点的顶点, 可以将其从 T 移开, 这并不会改变 T 的期望深度。这样, 剩下的是一棵树 T' , 有 $n!$ 片叶子, 每片叶子的概率为 $1/n!$ 。因为 $D(T') \geq n! \log n!$, 所以 T' (即 T) 的期望深度至少为 $(1/n!) n! \log n! = \log n!$ 。□

推论 每次比较排序平均至少需要进行 $c \log n$ 次比较, 其中 c 是常数。

有必要提及一种称为快速排序的有效算法。虽然其期望运行时间的下限是 $c n \log n$, 其中 c 是常数。但当运行在实际的机器上时, 它的期望运行时间与已知的其他基于比较的算法的运行时间成比例。快速排序算法有最坏情况的运行时间, 它是二次的, 但这对很多应用并不重要。

算法 3.5 快速排序。

输入: n 个元素的序列 S, a_1, a_2, \dots, a_n 。

输出: S 的有序排列元素。

方法: 使用图 3-7 中定义的递归过程 QUICKSORT。算法包括对 QUICKSORT(S) 的调用。□

定理 3.8 算法 3.5 在 $O(n \log n)$ 的期望时间内对 n 个元素的序列进行排序。

证明: 算法 3.5 的正确性可通过直接对 S 的规模进行归纳证明。为简单起见, 在时间分析中, 假设 S 的所有元素都是不同的。这个假设将会使在行 3 构建的 S_1 和 S_2 的规模最大化, 并因此使花费在行 4 的递归调用的平均时间最大化。设 $T(n)$ 是 QUICKSORT 排序 n 个元素的序列所需要的期望时间。显然, $T(0) = T(1) = b$, 其中 b 是常数。

```

procedure QUICKSORT(S):
1.  if S contains at most one element then return S
    else
        begin
2.          choose an element a randomly from S;
3.          let S1, S2, and S3 be the sequences of elements in S less
              than, equal to, and greater than a, respectively;
4.          return (QUICKSORT(S1) followed by S2 followed by
                  QUICKSORT(S3))
        end

```

图 3-7 快速排序程序

假设在行 2 中所选择的元素 a 是 n 个元素的序列 S 的第 i 小元素。那么, 行 4 的两个递归调用的期望时间分别是 $T(i-1)$ 和 $T(n-i)$ 。因为 i 在 1 和 n 之间的取值可能性是相同的, 而 QUICKSORT(S) 的平衡显然需要时间 cn , 其中 c 是一个常数。因此有关系:

$$T(n) \leq cn + \frac{1}{n} \sum_{i=1}^n [T(i-1) + T(n-i)], \text{ 对于 } n \geq 2 \quad (3-3)$$

对式(3-3)进行代数推导可得

$$T(n) \leq cn + \frac{2}{n} \sum_{i=0}^{n-1} T(i) \quad (3-4)$$

下面将证明, 当 $n \geq 2$ 时, $T(n) \leq kn \ln n$, 其中 $k = 2c + 2b$, 而且, $b = T(0) = T(1)$ 。对于归纳基础 $n = 2$, $T(2) \leq 2c + 2b$ 显然符合式(3-4)。对于归纳步, 将式(3-4)写为

$$T(n) \leq cn + \frac{4b}{n} + \frac{2}{n} \sum_{i=2}^{n-1} ki \ln i \quad (3-5)$$

因为 $i \ln i$ 是凹向上升的, 很容易地证明

$$\sum_{i=2}^{n-1} i \ln i \leq \int_2^n x \ln x \, dx \leq \frac{n^2 \ln n}{2} - \frac{n^2}{4} \quad (3-6)$$

用式(3-6)代替式(3-5), 可得

$$T(n) \leq cn + \frac{4b}{n} + kn \ln n - \frac{kn}{2} \quad (3-7)$$

因为 $n \geq 2$ 而且 $k = 2c + 2b$, 可以得出 $cn + 4b/n \leq kn/2$ 。这样, 根据式(3-7)可以得出 $T(n) \leq kn \ln n$ 。□

实际上还应该考虑两个细节。首先是在 QUICKSORT 的行 2“随机”选择元素 a 的方法。执行者可能为了简单总是选择序列 S 的第一个元素。这个选择会使 QUICKSORT 的性能比使用方程式(3-3)差很多。通常, 传给排序程序的序列是在“一定程度上”已经排好序的, 所以第一个元素的值比较小的概率高于平均概率。作为一个极端的例子, 读者可以验证, 如果 QUICKSORT 对已经排序好的不包含相同元素的序列进行排序时, 假设行 2 总是选择 S 的第一个元素 a , 那么序列 S_3 包含的元素总是比 S 少一个。在这种情况下, QUICKSORT 的步数代价是二次的。

在行 2 选择元素 a 的更好的方法是使用一个随机元素产生器产生一个整数 i , $1 \leq i \leq |S|$ [⊙], 然后, 在 S 中选择第 i 个元素作为划分元素 a 。简单些的方法是从 S 中选择一个元素样本, 然后使用元素样本的中位数作为划分元素。例如, 选择 S 的第一个元素、中间元素以及最后一个元素的中位数作为划分元素。

第二个细节是怎样有效地将 S 划分为 S_1 、 S_2 和 S_3 三个序列。使数组 A 拥有所有 n 个初始元素是可能的也是所期望的。因为 QUICKSORT 递归地调用自身, 它的变量 S 总是在连续的数组项中, 可假定为 $A[f]$, $A[f+1]$, \dots , $A[l]$, $1 \leq f \leq l \leq n$ 。选择“随机”元素 a 后, 可以适当地划分

⊙ 用 $|S|$ 表示序列 S 的长度。

S 。也就是说, 能够将 S_1 移到 $A[f], A[f+1], \dots, A[k]$; $S_2 \cup S_3$ 移到 $A[k+1], A[k+2], \dots, A[l]$, 取某个 $k(f \leq k \leq l)$ 。那么, 如果需要, $S_2 \cup S_3$ 也能够分裂, 但是除非 S_1 和 $S_2 \cup S_3$ 中的一个为空, 否则对于 S_1 和 $S_2 \cup S_3$, 简单地递归调用 QUICKSORT 更有效。

也许适当划分 S 的最简单方法是用两个指针指向数组 i 和 j 。开始时, $i=f$, 从 $A[f]$ 到 $A[i-1]$ 始终包含 S_1 的元素。类似地, 最初 $j=l$, 从 $A[j+1]$ 到 $A[l]$ 始终包含着 $S_2 \cup S_3$ 的元素。图 3-8 是执行划分的程序。

```

begin
1.   i ← f;
2.   j ← l;
3.   while i ≤ j do
      begin
4.         while A[j] ≥ a and j ≥ f do j ← j - 1;
5.         while A[i] < a and i ≤ l do i ← i + 1;
6.         if i < j then
            begin
7.                 interchange A[i] and A[j];
8.                 i ← i + 1;
9.                 j ← j - 1
            end
      end
end

```

图 3-8 适当划分 S 为 S_1 和 $S_2 \cup S_3$

划分后, 对数组 $A[f]$ 到 $A[i-1]$ 即 S_1 和数组 $A[j+1]$ 到 $A[l]$ 即 $S_2 \cup S_3$, 递归调用 QUICKSORT。然而, 如果 $i=f$ 即 $S_1 = \emptyset$, 必须首先从 $S_2 \cup S_3$ 移除至少一个 a 的实例, 而移除按其进行划分的元素是很方便的。应当注意, 如果这个数组用来表示序列, 就能够通过简单地传递指针给 QUICKSORT 作为参数, 其中指针指向所使用的部分数组的第一个和最后的存储位置。

例 3.5 用元素 $a=3$ 划分数组 A

1	2	3	4	5	6	7	8	9
6	9	3	1	2	7	1	8	3

行 4 的 **while** 语句使 j 从初值 9 下降为 7, 因为 $A[9]=3$ 和 $A[8]=8$ 都大于或等于 a , 但 $A[7]=1 < a$ 。 i 的初值为 1, 因为 $A[1]=6 \geq a$, 所以行 5 没有从初值 1 增加 i 的值。因此, 交换 $A[1]$ 和 $A[7]$, 令 i 值为 2, j 值为 6, 这样留下图 3-9a 所示的数组。执行 3~9 行循环两次之后的结果如图 3-9b 和图 3-9c 所示。此时 $i > j$, 行 3 的 **while** 语句执行完毕。□

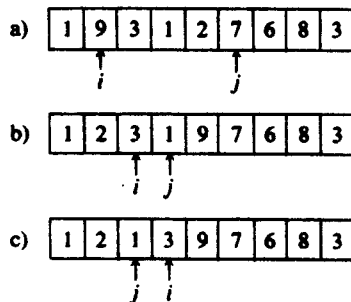


图 3-9 数组的划分

3.6 顺序统计学

和排序密切相关的一个问题是在 n 个元素的序列中选择第 k 小的元素。[⊖] 一个显然的解决方法

⊖ 严格地讲, 序列 a_1, a_2, \dots, a_n 的第 k 小的元素是序列中的元素 b , 有最多 $k-1$ 个 i 值使 $a_i < b$, 并且至少有 k 个 i 值使 $a_i \leq b$ 。例如, 4 是序列 7, 4, 2, 4 的第二和第三小的元素。

是排序序列为非递减序, 然后定位第 k 个元素。正如我们看到的, 这需要 $n \log n$ 次比较。通过细心应用分治策略, 可以用 $O(n)$ 步找到第 k 小的元素。一个重要的特殊例子是, 当 $k = \lceil n/2 \rceil$ 时, 可以在线性时间内找到一个序列的中位数。

算法 3.6 查找第 k 小元素。

输入: 从一个线性序的集合中所抽取的 n 个元素的序列和一个整数 k , $1 \leq k \leq n$ 。

输出: S 中第 k 小的元素。

方法: 使用图 3-10 的递归过程 SELECT。

□

```

1.  procedure SELECT( $k, S$ ):
2.      if  $|S| < 50$  then
3.          begin
4.              sort  $S$ ;
5.              return  $k$ th smallest element in  $S$ 
6.          end
7.      else
8.          begin
9.              divide  $S$  into  $\lfloor |S|/5 \rfloor$  sequences of 5 elements each
10.             with up to four leftover elements;
11.             sort each 5-element sequence;
12.             let  $M$  be the sequence of medians of the 5-element sets;
13.              $m \leftarrow \text{SELECT}(\lceil |M|/2 \rceil, M)$ ;
14.             let  $S_1, S_2$ , and  $S_3$  be the sequences of elements in  $S$  less
15.             than, equal to, and greater than  $m$ , respectively;
16.             if  $|S_1| \geq k$  then return SELECT( $k, S_1$ )
17.             else
18.                 if  $|S_1| + |S_2| \geq k$  then return  $m$ 
19.                 else return SELECT( $k - |S_1| - |S_2|, S_3$ )
20.             end
21.          end
22.      end

```

图 3-10 选择第 k 小元素的算法

直观地检验一下算法 3.6, 看它为什么有效。基本的想法是根据某个元素 m 将给定序列划分为三个序列 S_1 、 S_2 和 S_3 , 使得 S_1 包含所有小于 m 的元素, S_2 包含所有等于 m 的元素, S_3 包含所有大于 m 的元素。通过计算 S_1 和 S_2 中的元素数量, 能够确定第 k 小的元素在 S_1 、 S_2 还是 S_3 中。这样, 可以用一个小一点的问题代替给定的问题。

为了得到一个线性的算法, 必须能够在线性时间内找到划分元素, 使子序列 S_1 和 S_3 的大小不会大于 S 的大小的固定比例值。技巧就在于怎样选择划分元素 m 。将序列 S 划分为每个有 5 个元素的子序列。排序每个子序列, 用每个子序列的中位数组成序列 M 。现在 M 仅包含 $\lfloor n/5 \rfloor$ 个元素, 可以找到它的中位数, 比在 n 个元素序列中执行该操作要快 5 倍。

此外, S 中至少四分之一的元素小于或者等于 m ; 至少四分之一的元素大于或者等于 m 。这可以由图 3-11 说明。问题出现了, 为什么是“神奇数字”5? 答案是有两个 SELECT 的递归调用, 每一个都基于大小为 S 的一个比值的序列。为了使算法运行时间为线性, 两个序列的长度和必须小于 $|S|$ 。除了 5, 其他数也可以, 但是, 对某些数, 排序子序列的代价很高。确定哪些数适合于取代 5 留作习题。

定理 3.9 算法 3.6 在 $O(n)$ 时间内找到 n 个元素的序列 S 的第 k 小元素。

证明: 算法的正确性可对 S 的大小直接进行归纳证明, 这部分的证明留作练习。设 $T(n)$ 为从大小为 n 的序列中选择第 k 小元素需要的时间。中位数序列 M 的大小最多为 $n/5$, 这样, 递归调用

SELECT($\lceil |M|/2 \rceil, M$)

至多需要 $T(n/5)$ 时间。

由假设, 对于 R_p^+ , a_u 与 S_3 中的任何元素无关。假设在这个线性序 R 中, a_u 在 a_m 之前, 即 $a_u R a_m$ 。那么可以找到一个线性序 R' , 除了将 a_u 移到 a_m 后面, 作为 a_m 的直接后继, 它与 R 相同。 R' 也与 R_p^+ 一致。对每个 R 和 R' , 可以发现 a 的不同的整数值以分别满足 R 或者 R' 。但是, a_m 不能是两种情况下的第 k 小元素, 因为 a_m 在 R' 中比在 R 中的位置靠前一个元素。这样可以得出结论, 如果对于 R_p^+ , S 中的某个元素与 a_m 无关, 那么 T 没有在集合 S 中正确地选择第 k 个元素。对于情况 $a_m R a_u$, 可以对称地处理。□

定理 3.10 如果 T 是一棵在集合 S 中选择第 k 小元素的决策树, $\|S\| = n$, 那么 T 的每片叶子的深度至少为 $n-1$ 。

证明: 考虑 T 中从根到叶子的路径 p 。通过引理 3.3, 对每个 $i \neq m$, $a_i R_p^+ a_m$ 或者 $a_m R_p^+ a_i$, 其中将 a_m 选为第 k 小元素。对元素 a_i , $i \neq m$, 对 a_i 定义关键比较为 p 上包括 a_i 的第一次比较, 满足下述之一:

1. a_i 与 a_m 比较;
2. a_i 与 a_j 比较, $a_i R_p a_j$ 和 $a_j R_p^+ a_m$;
3. a_i 与 a_j 比较, $a_j R_p a_i$ 和 $a_m R_p^+ a_j$ 。

从直观上看, a_i 的关键比较是第一次比较, 根据它能够最终确定 a_i 在 a_m 之前或者之后。

显然, 除了 a_m , 每个元素 a_i 都有关键比较, 否则既没有 $a_i R_p^+ a_m$ 也没有 $a_m R_p^+ a_i$ 。而且容易看出两个被比较元素间的比较不会是关键比较。因为必须有 $n-1$ 个元素包含在关键比较中, 所以路径 p 的长度必须至少为 $n-1$ 。□

引理 在 S 中找到第 k 小元素, 不管在期望情况下还是最坏情况下, 都至少需要 $n-1$ 次比较。

实际上, 对所有 k , 除了 1 或 n , 可以证明比定理 3.10 更强的结果。见习题 3.21 ~ 3.23。

要在 S 中寻找计算第 k 小元素的好的期望时间的算法, 类似快速排序的策略是适用的。

算法 3.7 查找第 k 小元素。

输入: 从线性序 \leq 中抽出的 n 个元素的序列 S 和整数 k , $1 \leq k \leq n$ 。

输出: S 中的第 k 小元素。

方法: 使用图 3-12 所示的递归过程 SELECT。

```

procedure SELECT(k, S):
1.  if |S| = 1 then return the single element in S
    else
2.      begin
3.          choose an element a randomly from S;
4.          let S1, S2, and S3 be the sequences of elements in S less
              than, equal to, and greater than a, respectively;
5.          if |S1| ≥ k then return SELECT(k, S1)
              else
6.              if |S1| + |S2| ≥ k then return a
                  else return SELECT(k - |S1| - |S2|, S3)
              end
    end

```

图 3-12 选择算法

定理 3.11 算法 3.7 具有线性的期望运行时间。

证明: 设 $T(n)$ 是对 n 个元素序列执行 SELECT 的期望运行时间。为了简化, 假设 S 中所有元素都不同。(如果有重复元素, 结果并不会改变。)

假设行 2 选择的元素 a 是 S 中第 i 小的元素, 那么 i 可能与 $1, 2, \dots, n$ 中任意一个元素的概率是相同的。如果 $i > k$, 则对 $i-1$ 个元素的序列调用 SELECT; 如果 $i < k$, 则对 $n-i$ 个元素的

序列调用 SELECT。这样，行 4 或者行 6 的递归期望代价是

$$\frac{1}{n} \left[\sum_{i=1}^{k-1} T(n-i) + \sum_{i=k+1}^n T(i-1) \right] = \frac{1}{n} \left[\sum_{i=n-k+1}^{n-1} T(i) + \sum_{i=k}^{n-1} T(i) \right]$$

过程 SELECT 剩下的部分需要时间 cn ，其中 c 为常数，所以，对 $n \geq 2$ ，有以下不等式

$$T(n) \leq cn + \max_k \left\{ \frac{1}{n} \left[\sum_{i=n-k+1}^{n-1} T(i) + \sum_{i=k}^{n-1} T(i) \right] \right\} \quad (3-9)$$

读者可使用归纳法作为练习自行证明：如果 $T(1) \leq c$ ，那么对所有 $n \geq 2$ ， $T(n) \leq 4cn$ 。□

习题

- 3.1 用算法 3.1 排序字符串 $abc, acb, bca, bbc, acc, bac, baa$ 。
- 3.2 用算法 3.2 排序字符串 $a, bc, aab, baca, cbc, cc$ 。
- 3.3 根据例 3.2 测试图 3-13 中的两棵树是否同构。

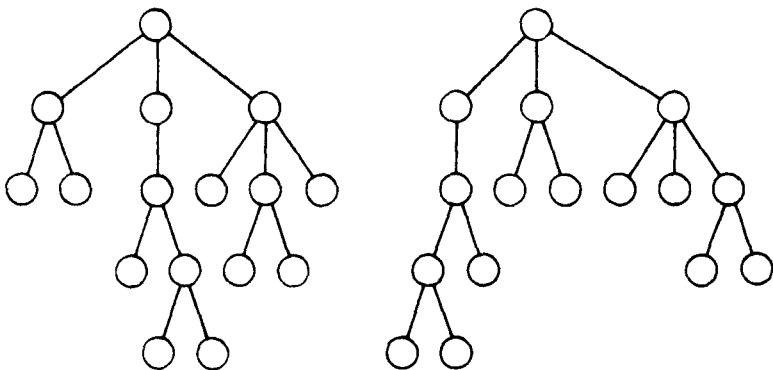


图 3-13 两棵树

- 3.4 排序表 3, 1, 4, 1, 5, 9, 2, 6, 5, 3, 5, 8, 9, 7，请分别使用 (a) 堆排序 (b) 快速排序和 (c) 归并排序 (算法 2.4)。每种情况需要比较几次？
- 3.5 考虑以下排序存储在数组 A 中的元素序列 a_1, a_2, \dots, a_n 的算法。也就是说， $A[i] = a_i$ ， $(1 \leq i \leq n)$ 。

```

procedure BUBBLESORT( $A$ ):
  for  $j = n - 1$  step  $-1$  until  $1$  do
    for  $i = 1$  step  $1$  until  $j$  do
      if  $A[i + 1] < A[i]$  then interchange  $A[i]$  and  $A[i + 1]$ 
  
```

- a) 证明 BUBBLESORT 按非递减次序对 A 中的元素进行排序。
- b) 确定 BUBBLESORT 的最坏情况运行时间和期望运行时间。
- 3.6 完成能够在线性时间内进行建堆的证明。(定理 3.5)
- 3.7 完成堆排序需要 $O(n \log n)$ 的时间的证明。(定理 3.6)
- 3.8 说明快速排序在最坏情况下的运行时间是 $O(n^2)$ 。
- 3.9 找第 k 小元素的算法 3.7 在最坏情况下的运行时间是什么？
- *3.10 通过完成定理 3.7 的证明以及解方程 (3-2)，证明排序 n 个元素的期望时间下界是 $cn \log n$ ，其中 c 为常数。
- 3.11 通过解方程 (3-8) 完成算法 3.6 在 $O(n)$ 时间内找到第 k 小元素 (定理 3.9) 的证明。
- *3.12 证明图 3-8 分割程序的正确性，并且分析它的运行时间。

- 3.13 通过解方程(3-9)证明找第 k 小元素的算法 3.7 的期望时间是 $O(n)$ (定理 3.11)。
- *3.14 说明算法 3.7 用的比较次数的期望值最多是 $4n$ 。如果知道算法用到的 k 值, 能否改进这个界限?
- *3.15 设 S 是元素序列, 其中第 i 个元素有 m_i 个副本, $1 \leq i \leq k$ 。设 $n = \sum_{i=1}^k m_i$, 证明若用比较排序对 S 进行排序,

$$O\left(n + \log\left(\frac{n!}{m_1! m_2! \cdots m_k!}\right)\right)$$

次比较是必要和充分的。

- 3.16 设 S_1, S_2, \dots, S_n 是 $1 \sim n$ 范围内的整数集合, 所有 S_i 的势的和是 n 。描述一个排序所有 S_i 的 $O(n)$ 的算法。
- 3.17 给定序列 a_1, a_2, \dots, a_n 和排列 $\pi(1), \pi(2), \dots, \pi(n)$, 编写简化 ALGOL 算法, 重新原地安排序列, 使得顺序为 $a_{\pi(1)}, a_{\pi(2)}, \dots, a_{\pi(n)}$ 。算法的最坏情况运行时间和期望运行时间是多少?
- *3.18 在构建大小为 $2^k - 1$ 的堆时, 可构建两个大小为 $2^{k-1} - 1$ 的堆, 再组合起来, 添加根并将根处的元素由根向下放到适当的位置。由此, 可以简单地通过一次增加一个元素作为新叶子, 将新元素向上推以添加到树上的方式来构建一个堆。写一个通过一次添加一个叶子构建堆的算法, 并与算法 3.3 比较渐近增长率。
- 3.19 考虑矩形数组。将每行的元素按升序排序。然后将每列元素按升序排序。证明每行元素仍然有序。
- *3.20 设 a_1, a_2, \dots, a_n 是一个元素序列, 设 p, q 是正整数。考虑选择每隔 p 的元素形成的子序列。排序这些子序列。对 q 重复这一过程。证明距离为 p 的子序列仍然有序。
- **3.21 考虑用比较的方法从 n 个元素的集合中查找最大和次大元素。证明 $n + \lceil \log n \rceil - 2$ 次比较是必要的和充分的。
- **3.22 证明在 n 个元素的序列中查找第 k 小元素所需的期望比较次数至少为 $(1 + 0.75\alpha(1 - \alpha))n$, 其中 $\alpha = k/n$, 其中 k 和 n 足够大。
- **3.23 说明在 n 个元素的集合中找到第 k 小元素, 在最坏情况下 $n + \text{MIN}(k, n - k + 1) - 2$ 次比较是必需的。
- **3.24 设 S 是 n 个整数的集合。假设只能执行 S 中元素相加以及和的比较。在这种情况下, 找到 S 的最大元素需要多少次比较?
- *3.25 算法 3.6 将序列分割为规模为 5 的子序列。算法对其他规模如 3, 7 和 9 是否也可行? 选择使得比较总次数最小的子序列规模。图 3-14 指出排序不同规模集合的已知最少比较次数。对 $n \leq 12$, 图中所示的比较次数是已知最优的。

规模 n	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17
比较次数	0	1	3	5	7	10	13	16	19	22	26	30	34	38	42	46	50

图 3-14 排序 n 个元素的已知最少比较次数

- *3.26 算法 3.5 将序列分为长度为 5 的子序列, 找到每个子序列的中位数, 并找到这些中位数的中位数。除了找到中位数的中位数, 寻找某个其他元素, 如第 $\lfloor k/5 \rfloor$ 个元素是否更高效?
- *3.27 考虑以下排序方法。以包含一个元素的序列开始, 通过折半查找一次插入剩下的一个元素。设计一个可以快速进行折半查找和插入元素的数据结构。是否可以在 $O(n \log n)$ 时间

内实现这种排序?

- ** 3.28 除了选择一个随机元素分割规模为 n 的集合外,如快速排序或者算法 3.7 所做的,还可以选择一个规模为 s 的小样本,找到它的中位数,利用中位数来分割整个集合。说明如何选择 s 作为 n 的函数,使排序所需的期望比较次数最小。
- ** 3.29 扩展习题 3.28 的思想,使顺序统计问题所需要的比较次数最小。[提示:从样本集合中选择两个元素,这两个元素具有较高的跨越所需元素的概率。]
- 3.30 如果相等元素排序后的顺序与排序前相同,则排序方法是稳定的。下列哪些排序算法是稳定的?
 - a) 冒泡排序(习题 3.5)
 - b) 归并排序(算法 2.4)
 - c) 堆排序(算法 3.4)
 - d) 快速排序(算法 3.5)

研究性问题

- 3.31 在特定情况下所需的比较次数有一些开放问题。例如,希望从 n 个元素的集合中查找第 k 小元素。Pratt 和 Yao[1973]已经讨论过 $k=3$ 的情况。对于 n 个元素, $n \geq 13$, 图 3-14 给出的数值是否最优尚是未知的。对比较小的 n , Ford 和 Johnson[1959]提出的排序算法在比较次数上是最优的。

文献及注释

Knuth[1973a]概述了排序方法。堆排序由 Williams[1964]提出,由 Floyd[1964]改进。快速排序根据 Hoare[1962]而来。习题 3.28 中对快速排序的改进是由 Singleton[1969]以及 Frazer 和 McKellar[1970]提出的。算法 3.6 查找顺序统计的线性最坏情况算法是由 Blum, Floyd, Pratt, Rivest 和 Tarjan[1972]提出的。Hadian 和 Sobel[1969]以及 Pratt 和 Yao[1973]讨论了查找特定的顺序统计值的比较次数。

习题 3.21 的结果可参见 Kislitsyn[1964]。习题 3.22、3.23 和 3.29 来自 Floyd 和 Rivest[1973],该文献还包括了一个比习题 3.22 所阐述的更强的下限。习题 3.19、3.20 和一些推广在 Gale 和 Karp[1970]和 Liu[1972]中有所讨论。排序的一个有趣的应用是在平面上查找点集合的凸包,它由 Graham[1972]提出。Horvath[1974]则讨论了稳定算法。

第4章 集合操作问题的数据结构

对于给定问题,探讨有效算法设计的一个较好方式是研究问题的基本性质。通常,一个问题可以表述为集合等基本数学对象,该问题的算法可以依据基于这些基础对象的基本操作来描述。这种观点的优势在于,通过研究各种数据结构,选择在整体上最适合于问题的算法。因此,好的数据结构设计能设计好的算法。

本章研究集合的7个基本操作,这些操作是许多查找和信息检索等问题特有的。本章介绍表示集合的多种数据结构,并考虑在需要执行一系列不同类型的操作时,每种结构的适用性。

4.1 集合的基本操作

考虑下述集合的基本操作:

1. MEMBER(a, S), 判定 a 是否为集合 S 中的元素。若是, 输出“是”; 否则, 输出“否”。
2. INSERT(a, S), 用 $S \cup \{a\}$ 替换集合 S 。
3. DELETE(a, S), 用 $S - \{a\}$ 替换集合 S 。
4. UNION(S_1, S_2, S_3), 用 $S_3 = S_1 \cup S_2$ 替换集合 S_1 和 S_2 。为避免需要在 $S_1 \cup S_2$ 中删除相同的元素, 当执行该操作时, 假定 S_1 和 S_2 不相交。
5. FIND(a), 输出包含元素 a 的集合的名称。如果 a 不止在一个集合中出现, 则该指令是不确定的。
6. SPLIT(a, S), 假定集合 S 中的元素关于“ \leq ”是线性有序的。该操作将 S 分为两个集合 S_1 和 S_2 , 使得 $S_1 = \{b \mid b \leq a \text{ 且 } b \in S\}$, $S_2 = \{b \mid b > a \text{ 且 } b \in S\}$ 。
7. MIN(S), 输出集合 S 中(关于线性序“ \leq ”)的最小元素。

在实际操作中所碰到的许多问题都可以被简化为一个或多个子问题, 每个子问题抽象表示成在某个数据基(元素的全集)上的基本指令序列。本章将考虑指令系列 σ , σ 中的指令取自上述7种集合操作的某个子集。

例如, MEMBER、INSERT 和 DELETE 操作的处理序列是许多查找问题的主要部分。能够用来处理 MEMBER、INSERT 和 DELETE 操作序列的数据结构称为字典(dictionary)。本章将研究几类用来实现字典的数据结构, 诸如散列表(hash table)、二叉查找树(binary search tree)和2-3树(2-3tree)等。

有许多有趣的问题, 这里主要关注 σ 的时间复杂度, 也就是用函数度量执行 σ 中的指令所需要的时间, 该函数与 σ 的长度和数据基的大小有关。我们将考虑平均和最坏情况下的时间复杂度, 并进一步区分在线(on-line)和离线(off-line)复杂度的差异。

定义 σ 的在线执行要求 σ 中的指令从左向右执行, 在 σ 中执行第 i 个指令时无需考虑任何后续指令。 σ 的离线执行允许在需要得到答案之前扫描 σ 中的所有指令。

显然, 任何在线算法可以用作离线算法, 但反过来则不一定。我们将注意到, 离线算法比任何已知的在线算法快。然而在许多应用中, 仅限于考虑在线算法。

给定要执行的一串指令, 最根本的问题是使用什么数据结构来表示基本的数据基。通常, 需要细心地平衡两个相互冲突的要求。典型情况下, 指令序列是若干以未知顺序重复执行的操作。有可能存在几个数据结构, 每个都能使某个操作执行起来很容易, 但其他操作执行困难。在多数情况下, 最好的解决方案是取某种折衷。通常会采用一些数据结构, 虽然它们不能让任何一

种操作最容易,但是却使整体的执行优于任何先前的方法。

现在举例说明如何根据集合操作的指令系列来表示图的生成树问题。

定义 给定无向图 $G=(V, E)$ 。 G 的生成树(spanning tree)是无向树 $S=(V, T)$ 。 $G=(V, E)$ 的生成森林是无向树的集合 $\{(V_1, T_1), (V_2, T_2), \dots, (V_k, T_k)\}$, 所有 V_i 组成了 V 的一个划分^①, 并且每个 T_i 都是 E 的(可能为空的)子集。

图 $G=(V, E)$ 的代价函数 c 是从 E 到实数的映射, 图 G 的子图 $G'=(V', E')$ 的代价 $c(G')$ 为 $\sum_{e \in E'} c(e)$ 。

例 4.1 考虑图 4-1 中的算法, 寻找给定图 $G=(V, E)$ 的一棵最小代价生成树 $S=(V, T)$ 。5.1 节将详细讨论这一最小代价生成树算法, 例 5.1 举例说明其应用。

```

begin
1.    $T \leftarrow \emptyset$ ;
2.    $VS \leftarrow \emptyset$ ;
3.   for each vertex  $v \in V$  do add the singleton set  $\{v\}$  to  $VS$ ;
4.   while  $\|VS\| > 1$  do
       begin
5.         choose  $(v, w)$ , an edge in  $E$  of lowest cost;
6.         delete  $(v, w)$  from  $E$ ;
7.         if  $v$  and  $w$  are in different sets  $W_1$  and  $W_2$  in  $VS$  then
            begin
8.             replace  $W_1$  and  $W_2$  in  $VS$  by  $W_1 \cup W_2$ ;
9.             add  $(v, w)$  to  $T$ 
            end
       end
end
end

```

图 4-1 最小代价生成树算法

图 4-1 的算法使用了 E 、 T 和 VS 3 个集合。集合 E 包含给定图 G 的边, 集合 T 用来存放最终生成树的边。算法将 G 的生成森林转化为一棵生成树。集合 VS 包含生成森林中树的顶点集合。初始时, VS 包含由 G 中的每个顶点自身所组成的集合。

将算法看作是处理 3 个集合 E 、 T 和 VS 的操作序列。第 1 行初始化集合 T 。第 2、3 行初始化 VS , VS 中的元素本身就是集合。第 3 行将初始的单元元素集添加到 VS 。第 4 行控制算法的主循环, 对集合 VS 中顶点集合的数目进行计数。第 7 行判定边 (v, w) 是否连接生成森林中的两棵树。若是, 则在第 8 行合并这两棵树, 并且在第 9 行将边 (v, w) 添加到最终的生成树中。

第 7 行要求能够在 VS 中找到含有特定顶点的集合的名字。(在 VS 中, 集合所使用的实际名称并不重要, 因此可使用任意的集合名。)基本上, 在第 7 行要求能够有效处理 FIND 原语。类似地, 第 8 行要求能够对不相交的顶点集合执行 UNION 操作。

要为单独处理 UNION 操作或 FIND 操作寻找数据结构是相当容易的。然而, 这里需要的数据结构应该使 UNION 和 FIND 操作都易于实现。进一步, 由于第 8 行 UNION 操作的执行依赖于第 7 行 FIND 操作的输出, 所需的 UNION 和 FIND 指令序列必须是在线执行的。在 4.6 和 4.7 节将研究两个这种类型的结构。

考虑对边集 E 执行的操作序列, 在第 5 行需要 MIN 询问原语, 第 6 行则需要 DELETE 原语。我们已经为这两个原语找到了一个好的数据结构, 就是 3.4 节介绍的堆(虽然那里用堆来找出最大元素, 但是显然也很容易用堆来找出最小元素)。

最后, 仅在第 9 行生成树的边集 T 要求 INSERT 操作, 增加一条新的边到 T 。用一个单链表即可满足。□

① 即 $V_1 \cup V_2 \cup \dots \cup V_k = V$, 并且对于 $i \neq j$, $V_i \cap V_j = \emptyset$ 。

4.2 散列法

在选择合适数据结构来处理指令序列 σ 时,除了考虑在给定指令序列 σ 中会出现哪些指令外,还需考虑的另一个重要问题是由 σ 所操作的数据基(全集)的大小。例如在第3章中,如果元素是某个适当范围内的整数,用桶排序可以在线性时间内完成 n 个元素序列的排序问题。然而,如果元素取自任意的线性有序集,那么可得到的最优时间为 $O(n \log n)$ 。

本节将考虑处于变化状态的元素集合 S 的维护问题。新元素添加到 S 中,从 S 中移除旧元素,并且不时地回答问题:“元素 x 当前在 S 中吗?”。很自然地,这个问题可由一个字典来模拟;我们需要一个数据结构,它允许便利地处理 MEMBER、INSERT 和 DELETE 指令系列。假定从一个非常大的全集中选择可能出现在 S 中的元素,将 S 表示为一个位向量就变得不切实际了。

例 4.2 编译程序或汇编程序会跟踪其正在翻译的程序中出现的所有标识符组成的“符号表”。对于大多数程序语言来说,所有可能标识符组成的集合是非常大的。例如,在 FORTRAN 语言中,存在大约 1.62×10^9 个可能的标识符。因此,用一个数组来表示符号表,并对每个可能标识符分配一项,而不考虑标识符是否确实出现在程序中,这是不可思议的。

对符号表,编译程序执行的操作分为两类。首先,当碰到新标识符时,必须将它们加入到表中。这包括在表中创建一个位置,存储特定标识符及其有关的数据(诸如实型还是整型?)。其次,编译程序随时可能请求关于某个标识符的信息(例如,标识符是整型吗?)。

因此,能够处理 INSERT 和 MEMBER 操作的数据结构可能必须足以实现一个符号表。实际上,本节讨论的数据结构经常用来实现一个符号表。□

考虑散列技术(hashing),不仅处理在符号表构造中所需要的 INSERT 和 MEMBER 指令,而且也处理 DELETE 指令。关于散列方案,存在多种变形,这里仅仅考虑其基本思想。

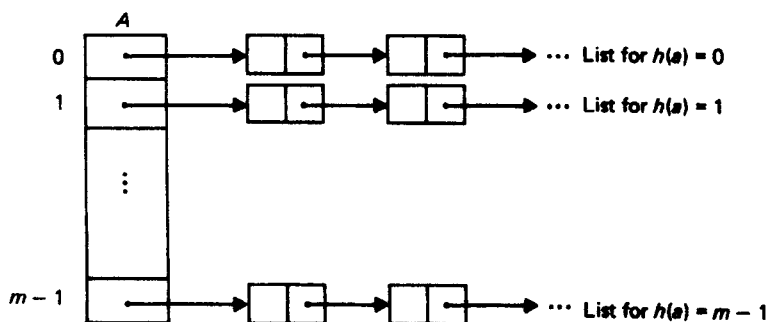


图 4-2 一种散列方案

图 4-2 给出了一种散列方案,其中的散列函数 h 将全集(例如,符号表中所有可能标识符的集合)中的元素映射到整数 0 到 $m-1$ 。假定对所有的元素 a ,能够在常数时间内计算 $h(a)$ 。存在大小为 m 的数组 A ,其项为指向 S 中元素列表的指针。由 $A[i]$ 指向的列表包含所有使得 $h(a) = i$ 的、在 S 中的元素 a 。

为执行指令 $\text{INSERT}(a, S)$,计算 $h(a)$ 并搜索由 $A[h(a)]$ 指向的列表。若 a 不在该列表中,则添加到表尾。为执行指令 $\text{DELETE}(a, S)$,再次搜索列表 $A[h(a)]$,如果 a 在列表中,则删除之。相似地, $\text{MEMBER}(a, S)$ 指令通过扫描列表 $A[h(a)]$ 来得到结果。

这种散列方案的计算复杂性是易于分析的。从最坏情况的角度,结果并不很好。例如,假定由 n 个不同的 INSERT 操作组成的序列 σ , h 应用到每个要插入的元素后可能产生相同的数,

结果所有元素都出现在同一列表中。在这种情况下, 处理 σ 中第 i 个指令需要与 i 成比例的时间。因此, 将所有 n 个元素加入到集合 S 中, 用散列法所需要的时间与 n^2 成比例。

然而, 从期望时间的角度看, 散列法看起来更好。若假定 $h(a)$ 等可能地为 $0 \sim m-1$ 中的任意值, 并且要插入 $n \leq m$ 个元素, 那么当插入第 i 个元素时, 被取代列表的期望长度为 $(i-1)/m$ (总是小于 1)。因此, 插入 n 个元素所需要的期望时间是 $O(n)$ 。如果执行 $O(n)$ 个 DELETE 和 MEMBER 操作连同 INSERT 操作, 总的期望代价仍然为 $O(n)$ 。

记住, 该分析假定散列表的规模 m 大于或等于集合 S 的最大规模 n 。但是, n 通常是预先未知的。当 n 未知时, 执行处理的合理方式是准备好构建一个散列表序列 T_0, T_1, T_2, \dots 。

选择合适的值作为初始散列表 T_0 的规模 m 。一旦插入到 T_0 的元素数目超过 m , 则创建一个规模为 $2m$ 的新散列表 T_1 , 并通过再次散列[⊖], 将当前处于 T_0 中的所有元素移入 T_1 。此时, 放弃旧的散列表 T_0 。接着在 T_1 继续插入更多的元素, 直到元素数目超过 $2m$ 。这时, 创建规模为 $4m$ 的新散列表 T_2 , 并将 T_1 中的元素重散列到 T_2 。一般而言, 一旦表 T_{k-1} 包含 $2^{k-1}m$ 个元素, 就创建规模为 $2^k m$ 个元素的散列表 T_k 。依此类推, 直到已经插入所有元素。

使用这种方案并假定 $m=1$, 考虑插入 2^k 个元素到散列表中需要的期望时间。可通过递归来模拟这一过程:

$$T(1) = 1$$

$$T(2^k) = T(2^{k-1}) + 2^k$$

显然, 其解为 $T(2^k) = 2^{k+1} - 1$ 。

由此可见, 由 n 个 INSERT、MEMBER 和 DELETE 指令组成的指令序列, 用散列法来处理, 可以在 $O(n)$ 期望时间内完成。

散列函数 h 的选择相当重要。如果要添加到 S 中的元素是均匀分布在 0 到 r 范围内的整数, 其中 $r > n$, 那么 $h(a)$ 可取模 m 余 a , 其中 m 为当前散列表的规模。散列函数的其他例子可在本章结尾所引用的一些参考文献中找到。

4.3 二分搜索

本节将讨论一个简单的搜索问题, 并比较 3 种不同解决方案的优劣。给定集合 S , 包含取自某个较大的全集的 n 个元素, 现在来处理仅包含 MEMBER 指令的操作序列 σ 。

最直观的解决方案是将 S 中的元素存储在一个列表中, 通过顺序搜索该列表来处理每个 MEMBER(a, S) 指令, 直到找到给定元素 a , 或者列表中的所有元素都已经查验为止。在最坏情况和期望情况下, 用这种解决方案处理 σ 中的所有指令, 所需要的时间与 $n \times |\sigma|$ 成比例。该方案的主要优点在于, 所需要的处理时间很少。

另一种解决方案是将 S 中的元素插入到规模为 $\|S\|$ 的散列表中, 通过搜索列表 $h(a)$ 来执行指令 MEMBER(a, S)。如果能找到一个好的散列函数 h , 则处理序列 σ 时, 该解决方案要求 $O(|\sigma|)$ 的期望时间和 $O(n|\sigma|)$ 的最坏情况时间。主要的难点在于找到一个散列函数, 使 S 中的元素均匀分布在整个散列表中。

如果 S 存在一个线性序 \leq , 则第 3 种解决方案是采用二分搜索。将 S 中的元素存储在数组 A 中, 对该数组进行排序, 使得 $A[1] < A[2] < \dots < A[n]$ 。然后判定是否元素 a 在 S 中。比较 a 和存储在位置 $\lfloor (1+n)/2 \rfloor$ 上的元素 b 。如果 $a=b$, 比较终止并回答“是”; 否则, 如果 $a < b$, 对数组的前半部分重复这一过程; 如果 $a > b$, 则对数组的后半部重复这个过程。通过不断地折半分离搜索域, 找到 a 或者判定其不在 S 中永远不会需要超过 $\lceil \log(n+1) \rceil$ 次比较。

⊖ 必须使用一个新的散列函数, 其散列值从 0 到 $2m-1$ 。

图 4-3 中给出的递归程序 $\text{SEARCH}(a, f, l)$ 在数组 A 的位置 $f, f+1, f+2, \dots, l$ 上查找元素 a 。为了确定 a 是否在 S 中, 调用 $\text{SEARCH}(a, 1, n)$ 。

要理解该程序如何工作, 可以设想将数组 A 表示为一棵二叉树。树的根位于位置 $\lfloor (1+n)/2 \rfloor$, 根的左右儿子位于位置 $\lfloor (1+n)/4 \rfloor$ 和 $\lfloor 3(1+n)/4 \rfloor$, 依此类推。下一节将给出二搜索的更为清晰的解释。

显然, 当查找 A 中任意元素时, SEARCH 至多包含 $\lceil \log(n+1) \rceil$ 次比较, 由于在树中不存在长于 $\lceil \log(n+1) \rceil$ 的路径。如果所有元素成为某次搜索的目标的可能性是相等的, 那么也可以证明(习题 4.4), 在处理序列 σ 中的 MEMBER 指令时, 过程 SEARCH 给出了最优的期望比较次数(也就是, $|\sigma| \times \log n$)。^①

```

procedure SEARCH(a, f, l):
if f > l then return "no"
else
  if a = A[⌊(f+l)/2⌋] then return "yes"
  else
    if a < A[⌊(f+l)/2⌋] then
      return SEARCH(a, f, ⌊(f+l)/2⌋ - 1)
    else return SEARCH(a, ⌊(f+l)/2⌋ + 1, l)

```

图 4-3 二分搜索算法

4.4 二叉查找树

考虑如下问题。要向集合 S 中插入元素, 并从 S 中删除元素。另外, 有时要知道一个给定的元素是否在 S 中, 或者当前在 S 中的最小元素是什么。假定要增加到 S 中的元素来自一个大的全集, 该全集对关系“ \leq ”是线性有序的。该问题能够抽象为处理 INSERT 、 DELETE 、 MEMBER 和 MIN 的指令序列。

我们已经知道, 对于处理由 INSERT 、 DELETE 和 MEMBER 组成的指令序列, 散列表是一个很好的数据结构。然而, 不搜索整个表, 就不可能在散列表中找到最小元素。适合所有四种指令的数据结构是二叉查找树。

定义 集合 S 的二叉查找树(binary search tree)是一棵标记二叉树, 用元素 $l(v) \in S$ 来标记每个顶点 v , 使得

1. 对 v 的左子树中的每个顶点 v , $l(u) < l(v)$,
2. 对 v 的右子树中的每个顶点 v , $l(u) > l(v)$, 并且
3. 对于每个元素 $a \in S$, 恰好存在一个顶点 v 使得 $l(v) = a$ 。

注意, 条件 1 和 2 意味着树的标记是有序的。条件 3 由 1 和 2 得出。

例 4.3 对于 ALGOL 关键字 **begin**、**else**、**end**、**if** 和 **then**, 图 4.4 给出了一棵可能的二叉查找树。这里的线性序是字典序。□

要确定元素 a 是否包含在由一棵二叉查找树表示的集合 S 中, 可将 a 与根的标记进行比较。如果根的标记是 a , 那么显然 a 在 S 中。如果 a 小于根的标记, 则搜索根的左子树(如果存在的话)。如果 a 大于根的标记, 则搜索根的右子树。要是 a 在树中, 那么它最终被找出。否则, 当搜索到空的左或右子树的时候, 过程终止。

算法 4.1 搜索一棵二叉查找树。

输入: 集合 S 的一棵二叉查找树 T 和元素 a 。

输出: 如果 $a \in S$, 则为“是”; 否则为“否”。

方法: 如果 T 为空^②, 返回“否”。否则, 令 r 为 T 的根。然后, 算法进行唯一的程序调用 $\text{SEARCH}(a, r)$, SEARCH 是图 4-5 定义的

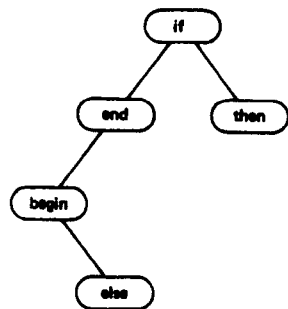


图 4-4 一棵二叉查找树

① 当然, 散列法的实现不仅通过比较, 而且可能散列法“优于”二分搜索, 并且在许多情况下, 确实如此。

② 虽然树的定义要求一棵树至少包含一个作为根的顶点, 但是在许多算法中, 将空树(不包含任何顶点的树)也作为二叉树来处理。

递归程序。

□

```

procedure SEARCH( $a, v$ ):
1. if  $a = l(v)$  then return "yes"
   else
2.   if  $a < l(v)$  then
3.     if  $v$  has a left son  $w$  then return SEARCH( $a, w$ )
4.     else return "no"
   else
5.     if  $v$  has a right son  $w$  then return SEARCH( $a, w$ )
6.     else return "no"

```

图 4-5 搜索一棵二叉查找树

显然, 执行指令 $\text{MEMBER}(a, S)$, 算法 4.1 已经足够了, 而且, 可以很容易地对它进行修改以执行指令 $\text{INSERT}(a, S)$ 。如果树为空, 则创建标记为 a 的根。如果树非空, 并且要插入的元素不存在树中, 那么在第 3 或 5 行, SEARCH 程序不能找到一个儿子。此时为元素创建一个新的顶点, 并关联到缺少儿子的位置, 而不是在第 4 或 6 行分别返回“否”。

执行 MIN 和 DELETE 指令时, 二叉查找树也很方便。沿着路径 v_0, v_1, \dots, v_p , 可找到二叉查找树 T 的最小元素, 其中 v_0 是 T 的根, 对于 $1 \leq i \leq p$, v_i 是 v_{i-1} 的左儿子, 且 v_p 没有左儿子, v_p 的标记是 T 中的最小元素。在某些问题中, 保存一个指向 v_p 的指针能提供便利, 以使得访问最小元素时仅需要常数访问时间。

指令 $\text{DELETE}(a, S)$ 的实现稍显困难。假定在顶点 v 找到了要删除的元素 a , 可能出现三种情形。

1. v 是叶子顶点, 在这种情况下, 从树中删除顶点 v ;
2. 顶点 v 恰好有一个儿子, 在这种情形, 使 v 的父亲成为 v 的儿子的父亲, 从树中删除 v (如果 v 是根顶点, 那么使 v 的儿子成为新的根);
3. 顶点 v 有两个儿子, 找到 v 的左子树中的最大元素 b 。递归地, 从子树中删除包含 b 的顶点, 并将 v 的标记设为 b 。注意, b 将成为整棵树中小于 a 的最大元素。

DELETE 操作的基于简化 ALGOL 语言的实现留作练习。注意, 单一的 MEMBER、INSERT、DELETE 或者 MIN 指令可能需用 $O(n)$ 时间。

例 4.4 假定要从图 4-4 的二叉查找树中删除单词 **if**。单词 **if** 位于根位置, 有两个儿子。在根的左子树中小于 **if** (字典序) 的最大单词是 **end**。从树中删除标记为 **end** 的顶点, 并且在根位置用 **end** 取代 **if**。然后, 由于 **end** 有一个儿子 (**begin**), 使 **begin** 成为根的儿子, 得到图 4-6 中的树。

□

当用一棵二叉查找树来表示基础集合时, 考虑由 n 个 INSERT 组成的指令序列的时间复杂度。将元素 a 插入到一棵二叉查找树中需要的时间是有界的, 它是 a 与树中已有元素的比较次数的常数倍。因此, 可以根据所作的比较次数来测量时间。

在最坏情形, 增加 n 个元素到一棵树中可能需要二次方时间。例如, 假定要增加的元素序列正巧是有序的 (增序)。在这种情况下, 查找树将包含由右儿子组成的单链。然而, 如果插入 n 个随机元素, 则像下述定理所表明的, 需要的插入时间为 $O(n \log n)$ 。

定理 4.1 对于 $n \geq 1$, 插入 n 个随机元素到初始为空的二叉查找树, 需要的期望比较次数是 $O(n \log n)$ 。

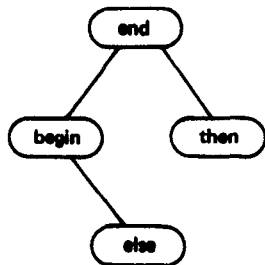


图 4-6 执行 DELETE 后的二叉查找树

证明：由序列 a_1, a_2, \dots, a_n 创建一棵二叉查找树时，令 $T(n)$ 为序列元素之间需要的比较次数。假定 $T(0) = 0$ ，并令 b_1, b_2, \dots, b_n 为已按升序排好的序列。

如果 a_1, \dots, a_n 为元素的随机序列，那么对于任意 $j (1 \leq j \leq n)$ ， a_j 等可能地为 b_j 。元素 a_j 成为二叉查找树的根，且在最终的树中， $j-1$ 个元素 b_1, b_2, \dots, b_{j-1} 将位于根的左子树， $n-j$ 个元素 $b_{j+1}, b_{j+2}, \dots, b_n$ 将包含在右子树中。

现在来计算将 b_1, b_2, \dots, b_{j-1} 插入到树中的期望比较次数。这些元素中的每个都同根比较一次，总共同根比较 $j-1$ 次。然后，递归地，插入 b_1, b_2, \dots, b_{j-1} 到树的左子树中，还需要 $T(j-1)$ 次比较。所有加起来，插入 b_1, b_2, \dots, b_{j-1} 到二叉查找树中，需要进行 $j-1 + T(j-1)$ 次比较。类似地，插入 $b_{j+1}, b_{j+2}, \dots, b_n$ 到树中，需要进行 $n-j + T(n-j)$ 次比较。

由于 j 等可能地取 1 到 n 中的任意值，有

$$T(n) = \frac{1}{n} \sum_{j=1}^n (n-1 + T(j-1) + T(n-j)) \quad (4-1)$$

对式(4-1)进行简单的代数运算，得到

$$T(n) = n-1 + \frac{2}{n} \sum_{j=0}^{n-1} T(j) \quad (4-2)$$

使用 3.5 节中的技术，能够证明

$$T(n) \leq k n \log n$$

其中 $k = \ln 4 = 1.39$ 。因此，插入 n 个元素到一棵二叉查找树中的期望比较次数是 $O(n \log n)$ 。□

总之，使用本节的技术，可以在 $O(n \log n)$ 的期望时间内，处理由 n 个 INSERT、DELETE、MEMBER 和 MIN 指令组成的随机指令序列。最坏情况下的执行性能是二次方的。然而，通过使用 4.9 节和习题 4.30~4.37 中讨论的平衡树方案(2-3 树、AVL 树或者有界平衡树)之一，即可使最坏情况下的性能提高到 $O(n \log n)$ 。

4.5 最优二叉查找树

在 4.3 节给定集合 $S = \{a_1, a_2, \dots, a_n\}$ (某大全集 U 的一个子集)，要求设计一个数据结构，可以有效地处理一个仅包含 MEMBER 指令的序列 σ 。再次考虑该问题，不过这次假定，除了给定的集合 S 外，对于全集 U 中所有的元素 a ，给出指令 MEMBER(a, S) 在 σ 中出现的概率。现在为 S 设计一棵二叉查找树，以期望的最小比较次数在线处理 MEMBER 指令序列 σ 。

假定 a_1, a_2, \dots, a_n 为集合 S 中的元素，以升序排列。 p_i 为指令 MEMBER(a_i, S) 在 σ 中出现的概率。对于某些 $a < a_1$ ， q_0 为形如 MEMBER(a, S) 的指令在 σ 中出现的概率。对于某些 $a_i < a < a_{i+1}$ ， q_i 为形如 MEMBER(a, S) 的指令在 σ 中出现的概率；对于某些 $a > a_n$ ， q_n 为形如 MEMBER(a, S) 的指令在 σ 中出现的概率。为了定义二叉查找树的代价(cost)，可加入 $n+1$ 个假想的叶子(fictitious leaves)，以便于构造集合 $U-S$ 中的元素的映射，称这些叶子为 0, 1, \dots, n 。

图 4-7 为图 4-4 的二叉查找树加入了这类假想叶子。例如，标记为 3 的叶子表示元素 a ，满足 $\text{end} < a < \text{if}$ 。

接下来需要定义二叉查找树的代价。如果元素 a 为某个顶点 v 的标记 $l(v)$ ，则在处理指令 MEMBER(a, S) 时访问的顶点数比顶点 v 的深度大 1。若 $a \notin S$ ，且 $a_i < a < a_{i+1}$ ，为处理指令 MEMBER(a, S)，访

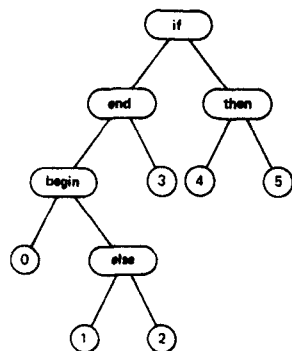


图 4-7 添加叶子后的
二叉查找树

问的顶点数等于假想叶子 i 的深度。因而, 二叉查找树的代价可如下定义

$$\sum_{i=1}^n p_i \times (\text{DEPTH}(a_i) + 1) + \sum_{i=0}^n q_i \times \text{DEPTH}(i)$$

一旦有了最小代价二叉查找树 T , 就能够以顶点访问的最小期望次数来执行 MEMBER 指令序列, 只需简单地将算法 4.1 应用于 T , 处理每条 MEMBER 指令。

给定各个 p_i 和 q_i , 如何找到一棵最小代价树呢? 分治 (divide-and-conquer) 法要求先确定根元素 a_i 。这就把问题分为两个子问题: 构造左子树和构造右子树。然而, 要确定根似乎并不容易, 因此, 解决整个问题还是有困难的。由于要考虑 $2n$ 个子问题, 而对应每个可能的根包含两个子问题, 由此, 很自然引出了动态规划方案。

对于 $0 \leq i < j \leq n$, T_{ij} 为元素 $\{a_{i+1}, a_{i+2}, \dots, a_j\}$ 的子集的一棵最小代价树。 c_{ij} 为 T_{ij} 的代价, r_{ij} 为 T_{ij} 的根。 T_{ij} 的权 (weight) w_{ij} 定义为 $q_i + (p_{i+1} + q_{i+1}) + \dots + (p_j + q_j)$ 。

如图 4-8 所示, 树 T_{ij} 包含根 a_k , 左子树 $T_{i,k-1}$ 为 $\{a_{i+1}, a_{i+2}, \dots, a_{k-1}\}$ 的一棵最小代价树, 右子树 $T_{k,j}$ 为 $\{a_{k+1}, a_{k+2}, \dots, a_j\}$ 的一棵最小代价树。若 $i = k - 1$, 则不存在左子树; 若 $k = j$, 无右子树。为便于符号表示, 将 T_{ii} 看作为空树, T_{ii} 的权值 w_{ii} 等于 q_i , 其代价 c_{ii} 为 0。

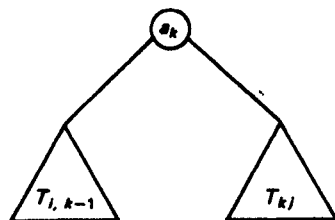


图 4-8 子树 T_{ij}

由树 $T_{i,k-1}$ 和 $T_{k,j}$ 中的深度, 树 T_{ij} 的左右子树上每个顶点的深度增 1, 其中 $i < j$ 。因此, T_{ij} 的代价 c_{ij} 可表示为

$$\begin{aligned} c_{ij} &= w_{i,k-1} + p_k + w_{k,j} + c_{i,k-1} + c_{k,j} \\ &= w_{ij} + c_{i,k-1} + c_{k,j} \end{aligned}$$

其中 k 的取值使得和 $c_{i,k-1} + c_{k,j}$ 为最小。因此, 为找到一棵最优树 T_{ij} , 对于每个 $k (i < k \leq j)$, 计算树的代价, 其中该树以 a_k 为根, 左子树为 $T_{i,k-1}$, 右子树为 $T_{k,j}$, 然后选择一棵具有最小代价的树。具体细节包含于下述算法中。

算法 4.2 构造一棵最优的二叉查找树。

输入: 元素集 $S = \{a_1, a_2, \dots, a_n\}$ 。假定 $a_1 < a_2 < \dots < a_n$, 给定概率 q_0, q_1, \dots, q_n 和 p_1, p_2, \dots, p_n , 使得对于 $1 \leq i < n$, q_i 表示执行一条 MEMBER(a, S) 指令的概率, 该指令接收所有满足 $a_i < a < a_{i+1}$ 的 a 。 q_0 表示 $a > a_n$ 时, 执行 MEMBER(a, S) 指令的概率。对于 $1 \leq i \leq n$, p_i 表示执行 MEMBER(a_i, S) 指令的概率。

输出: S 的一棵最小代价二叉查找树。

方法:

1. 对于 $0 \leq i < j \leq n$, 应用图 4-9 的动态规划算法, 按 $j - i$ 值的增序, 计算 r_{ij} 和 c_{ij} 。
2. 在计算了各个 r_{ij} 后, 调用 BUILDTREE(0, n) 为 T_{0n} 递归地构造一棵最优树。图 4-10 是程序 BUILDTREE。

□

例 4.5 考虑 4 个元素 $a_1 < a_2 < a_3 < a_4$, $q_0 = 1/8$, $q_1 = 3/16$, $q_2 = q_3 = q_4 = 1/16$, 且 $p_1 = 1/4$, $p_2 = 1/8$, $p_3 = p_4 = 1/16$ 。图 4-11 给出由图 4-9 所给出的算法计算出的 w_{ij} 、 r_{ij} 和 c_{ij} 的值。为便于表示, 此表中 w_{ij} 和 c_{ij} 的值都乘了 16。

比如, 要计算 r_{14} , 必须比较 $c_{11} + c_{24}$ 、 $c_{12} + c_{34}$ 和 $c_{13} + c_{44}$ 的值, (乘以 16 后) 分别为 8、9 和 11。因此, 在图 4-9 的第 8 行, $k = 2$ 时最小, 所以 $r_{14} = a_2$ 。

```

begin
1.   for  $i \leftarrow 0$  until  $n$  do
      begin
2.          $w_U \leftarrow q_i$ ;
3.          $c_U \leftarrow 0$ ;
      end;
4.   for  $l \leftarrow 1$  until  $n$  do
      for  $i \leftarrow 0$  until  $n - l$  do
      begin
6.          $j \leftarrow i + l$ ;
7.          $w_U \leftarrow w_{i,j-1} + p_j + q_j$ ;
8.         let  $m$  be a value of  $k$ ,  $i < k \leq j$ , for which  $c_{i,k-1} + c_{kj}$ 
           is minimum;
9.          $c_U \leftarrow w_U + c_{i,m-1} + c_{mj}$ ;
10.         $r_U \leftarrow a_m$ ;
      end
end

```

图 4-9 计算最优子树的根的算法

```

procedure BUILDTREE( $i, j$ ):
begin
  create vertex  $v_U$ , the root of  $T_U$ ;
  label  $v_U$  by  $r_U$ ;
  let  $m$  be the subscript of  $r_U$  (i.e.,  $r_U = a_m$ );
  if  $i < m - 1$  then make BUILDTREE( $i, m - 1$ ) the left subtree of  $v_U$ ;
  if  $m < j$  then make BUILDTREE( $m, j$ ) the right subtree of  $v_U$ ;
end

```

图 4-10 构造最优二叉查找树的程序

		$i \rightarrow$				
		0	1	2	3	4
$l = j - i$ \downarrow	0	$w_{00} = 2$ $c_{00} = 0$	$w_{11} = 3$ $c_{11} = 0$	$w_{22} = 1$ $c_{22} = 0$	$w_{33} = 1$ $c_{33} = 0$	$w_{44} = 1$ $c_{44} = 0$
	1	$w_{01} = 9$ $c_{01} = 9$ $r_{01} = a_1$	$w_{12} = 6$ $c_{12} = 6$ $r_{12} = a_2$	$w_{23} = 3$ $c_{23} = 3$ $r_{23} = a_3$	$w_{34} = 3$ $c_{34} = 3$ $r_{34} = a_4$	
	2	$w_{02} = 12$ $c_{02} = 18$ $r_{02} = a_1$	$w_{13} = 8$ $c_{13} = 11$ $r_{13} = a_2$	$w_{24} = 5$ $c_{24} = 8$ $r_{24} = a_4$		
	3	$w_{03} = 14$ $c_{03} = 25$ $r_{03} = a_1$	$w_{14} = 10$ $c_{14} = 18$ $r_{14} = a_2$			
	4	$w_{04} = 16$ $c_{04} = 33$ $r_{04} = a_2$				

图 4-11 w_{ij} 、 r_{ij} 和 c_{ij} 的值

计算了图 4-11 的表之后, 可以通过调用 BUILDTREE(0, 4) 构造树 T_{04} 。得到的二叉查找树如图 4-12 所示, 这棵树的代价为 33/16。□

定理 4.2 用算法 4.2 构造一棵最优二叉查找树, 需要 $O(n^3)$ 的时间。

证明: 计算了 r_{ij} 表, 就可以应用程序 BUILDTREE, 在 $O(n)$ 时间内从该表构造一棵最优树。程序仅有 n 次调用, 且每次调用需要的时间为常数。

代价最高的部分是图 4-9 中的动态规划算法。第 8 行寻找合适的 k 值, 使 $c_{i,k-1} + c_k$ 最小。这要求 $O(j-i)$ 时间。在第 5~10 行之间的循环中的其他步骤需要常数时间。第 4 行执行 n 次外部循环, 对于外循环的每次迭代, 至多执行 n 次内循环。因此, 总的代价是 $O(n^3)$ 。

就算法的正确性而言, 对 $l=j-i$ 的一个简单规约表明, 在第 9 和 10 行正确计算了 r_{ij} 和 c_{ij} 。

为了证明可通过 BUILDTREE 构造一棵最优树是正确的, 通过观察, 如果顶点 v_{ij} 是 $\{a_{i+1}, a_{i+2}, \dots, a_j\}$ 的一棵子树的根, 那么其左儿子将是 $\{a_{i+1}, a_{i+2}, \dots, a_{m-1}\}$ 的一棵最优树的根, 其中 $r_{ij} = a_m$, 且其右儿子将是 $\{a_{m+1}, a_{m+2}, \dots, a_j\}$ 的一棵最优树的根。因此, 显而易见, 由归纳证明, BUILDTREE(i, j) 为 $\{a_{i+1}, a_{i+2}, \dots, a_j\}$ 正确地构造了一棵最优树。□

在算法 4.2 中, 可以把图 4-9 第 8 行中对 m 的查找限定在 $T_{i,j-1}$ 和 $T_{i+1,j}$ 的根位置之间的范围, 且仍可保证找到一个最小值。通过这一修改, 算法 4.2 能够在 $O(n^2)$ 时间内找到一棵最优树。

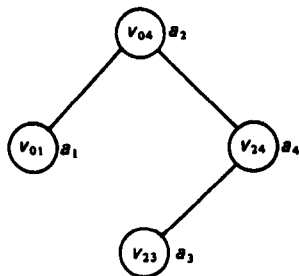


图 4-12 一棵最小代价树

4.6 简单的不相交集合并算法

考虑在例 4.1 的生成树算法中对顶点的处理, 由此引出的集合处理问题具有下列三个特征。

1. 任何时候合并两个集合, 它们都是不相交的;
2. 集合中的元素可以看作是 $1 \sim n$ 之间的整数;
3. 操作为 UNION 和 FIND。

本节和下一节将考虑这类问题的数据结构。假设给定 n 个元素, 取整数 $1, 2, \dots, n$ 。初始时, 假定每个元素各自处于本身的集合中。执行 UNION 和 FIND 指令序列。回想一下, UNION 指令的形式为 UNION(A, B, C), 表示两个不相交集 A 和 B 将被它们的并所替换, 该并集称为 C 。在应用中, 选用什么来命名集合通常无关紧要, 因此, 假定集合可以使用 $1 \sim n$ 之间的整数命名。而且, 假定不会给两个集合以相同的名字。

对于该问题的处理, 存在几种有趣的数据结构。本节将给出一个能够在 $O(n \log n)$ 时间处理指令序列的数据结构, 该指令序列至多可包含 $n-1$ 条 UNION 指令和 $O(n \log n)$ 条 FIND 指令。下一节将描述另一个数据结构, 它将在一个近乎线性于 n 的最坏情况时间内, 处理由 $O(n)$ 条 UNION 和 FIND 指令组成的序列。这些数据结构也能够以同样的计算复杂性, 处理 INSERT、DELETE 和 MEMBER 指令序列。

注意, 在 4.2~4.5 节中考虑的搜索算法假定元素取自一个全集, 该全集的规模远大于要执行指令的数目。本节假设全集的规模同要执行的指令序列的长度近似。

或许对于 UNION-FIND 问题, 最简单的数据结构是使用数组表示在任一时刻出现的集合的合集(collection of sets)。设 R 为大小为 n 的数组, $R[i]$ 表示包含元素 i 的集合名。由于集合的名字并不重要, 初始可以取 $R[i] = i, 1 \leq i \leq n$, 表示在初始时集合的合集可表示为 $\{1\}, \{2\}, \dots, \{n\}$, 并且集合 $\{i\}$ 命名为 i 。

执行指令 FIND(i) 可通过打印 $R[i]$ 的当前值实现。因此, 执行一条 FIND 指令的代价是常

量,这正是我们所期望的最好情形。

为执行指令 $\text{UNION}(A, B, C)$, 顺序扫描数组 R , 且包含 A 或 B 的每项都设置为 C 。因此, 执行一条 UNION 指令的代价是 $O(n)$ 。 n 个 UNION 指令序列可能需要 $O(n^2)$ 时间, 这种结果并不令人满意。

这一简单的算法可以通过几种方式加以改进, 一种改进是采用链表。另一种则是由于认识到将一个小的集合合并到更大的一个集合会更加有效。为此, 需要区分用来在数组 R 中标识集合的“内部名”和在 UNION 指令中提及的“外部名”。两者都可能是 $1 \sim n$ 之间的整数, 但并不一定相同。

针对该问题考虑下述数据结构。如前所述, 使用数组 R , 使 $R[i]$ 包含有元素 i 的集合的“内部”名。现在, 对于每个集合 A , 构造包含集合中元素的链表 $\text{LIST}[A]$ 。用两个数组 LIST 和 NEXT 实现该链表。 $\text{LIST}[A]$ 为整数 j , 表明 j 是以 A 为内部名的集合中的第一个元素。 $\text{NEXT}[j]$ 给出在集合 A 中的下一个元素, $\text{NEXT}[\text{NEXT}[j]]$ 表示再下一个元素, 依此类推。

此外, 使用数组 SIZE , $\text{SIZE}[A]$ 表示集合 A 中元素的数目。同样, 集合也可以在内部重命名。两个数组 INTERNAL_NAME 和 EXTERNAL_NAME 分别表示内部名和外部名。也就是, $\text{EXTERNAL_NAME}[A]$ 给出有内部名 A 的集合的真实名字(由 UNION 指令指定的名字)。 $\text{INTERNAL_NAME}[j]$ 是拥有外部名 j 的集合的内部名。内部名就是在数组 R 中使用的名字。

例 4.6 设 $n=8$, 且有合集, 包含三个集合 $\{1, 3, 5, 7\}$ 、 $\{2, 4, 8\}$ 和 $\{6\}$, 外部名分别为 1、2 和 3。图 4-13 给出这三个集合的数据结构, 其中, 假定 1、2 和 3 分别有内部名 2、3 和 1。□

R NEXT		
1	2	3
2	3	4
3	2	5
4	3	8
5	2	7
6	1	0
7	2	0
8	3	0

LIST	SIZE	EXTERNAL __NAME
6	1	3
1	4	1
2	3	2

集合 (有外部名)		INTERNAL __NAME
1 = {1, 3, 5, 7}	1	2
2 = {2, 4, 8}	2	3
3 = {6}	3	1

图 4-13 UNION-FIND 算法的数据结构

与先前一样执行指令 $\text{FIND}(i)$, 参照 $R[i]$ 来确定当前正包含元素 i 的集合的内部名。则 $\text{EXTERNAL_NAME}[R[i]]$ 给出包含 i 的集合的真实名。

如下执行形如 $\text{UNION}(I, J, K)$ 的合并指令(行号对应于图 4-14)。

1. 确定集合 I 和 J 的内部名(第 1~2 行)。
2. 参照数组 SIZE , 比较集合 I 和 J 的相对大小(第 3~4 行)。
3. 遍历较小集合的元素列表, 在数组 R 中改变相应项以指向较大集合的内部名(第 5~9 行)。
4. 把较小集元素表添加到较大集元素表的开始处, 将较小集合与较大集合(第 10~12 行)合并。

5. 将合并后集合的外部名命名为 K (第 13~14 行)。

通过将较小集合并到较大集,可以在与较小集合的元素个数成比例的时间内,执行 UNION 指令。图 4-14 中的程序给出了完整的细节。

```

procedure UNION( $I, J, K$ ):
begin
1.    $A \leftarrow \text{INTERNAL\_NAME}[I]$ ;
2.    $B \leftarrow \text{INTERNAL\_NAME}[J]$ ;
3.   wig assume  $\text{SIZE}[A] \leq \text{SIZE}[B]$ 
4.   otherwise interchange roles of  $A$  and  $B$  in
       begin
5.        $\text{ELEMENT} \leftarrow \text{LIST}[A]$ ;
6.       while  $\text{ELEMENT} \neq 0$  do
           begin
7.            $R[\text{ELEMENT}] \leftarrow B$ ;
8.            $\text{LAST} \leftarrow \text{ELEMENT}$ ;
9.            $\text{ELEMENT} \leftarrow \text{NEXT}[\text{ELEMENT}]$ 
           end;
10.       $\text{NEXT}[\text{LAST}] \leftarrow \text{LIST}[B]$ ;
11.       $\text{LIST}[B] \leftarrow \text{LIST}[A]$ ;
12.       $\text{SIZE}[B] \leftarrow \text{SIZE}[A] + \text{SIZE}[B]$ ;
13.       $\text{INTERNAL\_NAME}[K] \leftarrow B$ ;
14.       $\text{EXTERNAL\_NAME}[B] \leftarrow K$ 
       end
end

```

图 4-14 UNION 指令的实现

例 4.7 在执行指令 UNION(1, 2, 4)后,图 4-13 的数据结构将变为如图 4-15 所示的数据结构。□

	R	NEXT
1	2	3
2	2	4
3	2	5
4	2	8
5	2	7
6	1	0
7	2	0
8	2	1

	LIST	SIZE	EXTERNAL __NAME
1	6	1	3
2	2	7	4
3	-	-	-
4	-	-	-

集合 (有外部名)

3 = {6}

4 = {1, 2, 3, 4, 5, 7, 8}

	INTERNAL __NAME
1	-
2	-
3	1
4	2

图 4-15 执行 UNION 指令后的数据结构

定理 4.3 使用图 4-14 的算法,可以在 $O(n \log n)$ 步内执行 $n-1$ 个 UNION 操作(最大可能数)。

证明: 由于每个 UNION 的代价与移动的元素数成比例,在被移动的元素之间按比例均分每个 UNION 指令的代价,使得对于一个元素,每次移动都会有一个固定的代价计数。主要的观测

结果在于, 每次从列表中移动一个元素时, 它在一个至少两倍于之前长度的列表中找到自身。因此, 没有元素移动超过 $\log n$ 次, 因此, 计入任意一个元素的总代价是 $O(\log n)$ 。将计入到元素的代价相加, 就得到总代价 $O(n \log n)$ 。□

根据定理 4.3, 如果执行 m 条 FIND 指令和多达 $n-1$ 条 UNION 指令, 那么总的时间消耗为 $O(\text{MAX}(m, n \log n))$ 。如果 m 以 $n \log n$ 为阶或更大, 则该算法关于常数因子是最优的。然而, 在许多情况下可以发现 m 是 $O(n)$, 此时, 可以得到比 $O(\text{MAX}(m, n \log n))$ 更好的结果, 在下一节将能够看到这一点。

4.7 UNION-FIND 问题的树结构

对于 UNION-FIND 问题, 上节给出了一个数据结构, 能够在 $O(n \log n)$ 时间内处理 $n-1$ 条 UNION 指令和 $O(n \log n)$ 条 FIND 指令。将集合的合集元素表示成一棵树, 所有这些树将组成一个森林, 本节将给出表示森林的一个数据结构。该数据结构将允许在近乎线性的时间内, 处理 $O(n)$ 条 UNION 和 FIND 指令。

假定用一个有根的无向树 T_A 表示每个集合 A , A 中的元素对应于 T_A 中的顶点, 集合的名字同树的根联系。指令 $\text{UNION}(A, B, C)$ 的执行使 T_A 的根成为 T_B 的根的一个儿子, 并将 T_B 的根的名字改为 C 。指令 $\text{FIND}(i)$ 的执行在森林的某棵树 T 定位表示元素 i 的顶点, 遍历从该顶点到树 T 的根的路径, 由此可找到包含元素 i 的集合的名字。

使用这一方案, 合并两棵树的代价是常数。然而, $\text{FIND}(i)$ 指令的代价与从顶点 i 到根的路径长度同阶。这条路径的长度可能为 $n-1$ 。因此, 执行 $n-1$ 条 UNION 指令, 紧接着执行 n 条 FIND 指令, 代价可能高达 $O(n^2)$ 。例如, 考虑下面指令序列的代价:

```

UNION(1, 2, 2)
UNION(2, 3, 3)
...
UNION(n-1, n, n)
FIND(1)
FIND(2)
...
FIND(n)

```

执行 $n-1$ 条 UNION 指令得到图 4-16 中的树。 n 条 FIND 指令的代价正比于

$$\sum_{i=0}^{n-1} i = \frac{n(n-1)}{2}$$

如果能够保持树的平衡, 则可以减少代价。一种实现方式是, 计算每棵树中的顶点数, 在合并两个集合的时候, 总是将较小的树加到较大树的根上。这与上一节所使用的技术相似, 都是将较小集合合并到较大集。

引理 4.1 如果在执行每条 UNION 指令时, 顶点数较少(任意约定)的树的根成了较大树的根的儿子, 那么在森林中没有树的高度会大于等于 h , 除非其至少有 2^h 个顶点。

证明: 证明基于对 h 的归纳。对 $h=0$, 假设成立, 因为每棵树至少有一个顶点。假定对于所有小于 $h \geq 1$ 的值, 归纳假设成立。设 T 是高为 h 且顶点数最少的一棵树。则 T 必定是通过合并两棵树 T_1 和 T_2 得到的, 其中, T_1 的高为 $h-1$, 并且 T_1 的顶点数不多于 T_2 。由归纳假设, T_1 至

少有 2^{h-1} 个顶点, 因此 T_2 至少有 2^{h-1} 个顶点, 这表明 T 至少有 2^h 个顶点。

使用森林数据结构, 对于包含 n 条 UNION 和 FIND 的指令序列, 考虑其最坏情况下的执行时间。所作的修改是, 在一条 UNION 中, 较小树的根变成较大树的根的儿子。没有高度超过 $\log n$ 的树。因此, 执行 $O(n)$ 条 UNION 和 FIND 指令, 至多消耗 $O(n \log n)$ 个单元的时间。这个界是紧的, 因为存在包含有 n 条指令的序列, 使得消耗的时间正比于 $n \log n$ 。

现在考虑算法的另一种修改, 即所谓的路径压缩 (path compression)。由于在总代价中, FIND 指令的代价看起来占主导地位, 尝试减少 FIND 指令的代价。每次执行一条 FIND 指令, 都遍历由顶点 i 到对应树的根 r 的路径。设 $i, v_1, v_2, \dots, v_n, r$ 为该路径上的顶点。然后, 让 $i, v_1, v_2, \dots, v_{n-1}$ 中的每个都成为根的儿子。图 4-17b 给出将 FIND(i) 指令应用于图 4-17a 中的树之后的效果。

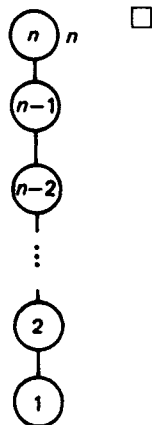


图4-16 执行UNION指令后的树

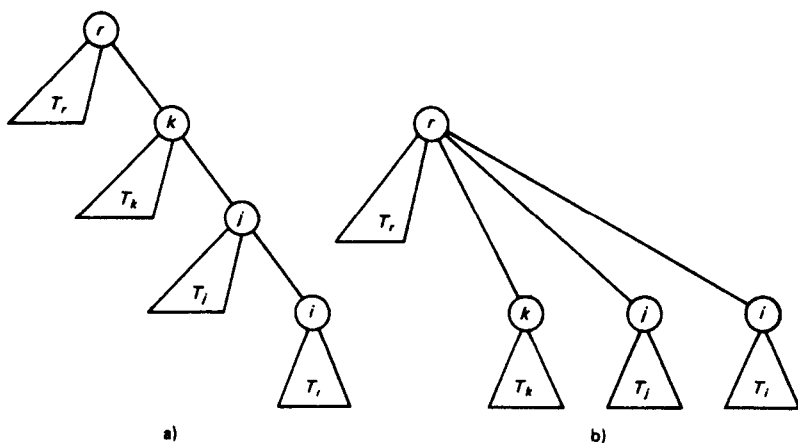


图 4-17 路径压缩的效果

对于 UNION-FIND 问题, 包含路径压缩的完整树合并 (tree-merging) 算法表述如下。

算法 4.3 不相交集的快速合并算法。

输入: 由集合的合集上的 UNION 和 FIND 指令组成的指令序列 σ , 集合中的元素通过取自 $1 \sim n$ 中的整数组成。集合的名字也假定为 $1 \sim n$ 中的整数, 初始时, 元素 i 本身处于名为 i 的集合中。

输出: 相应于 σ 中 FIND 指令的响应序列。在找到 σ 的下一条指令之前, 生成每条 FIND 指令的响应结果。

方法: 分 3 个部分来描述算法: 初始化、FIND 的响应和 UNION 的响应。

1. **初始化。** 对每个元素 $i (1 \leq i \leq n)$, 创建一个顶点 v_i 。设定 $\text{COUNT}[v_i] = 1$, $\text{NAME}[v_i] = i$ 和 $\text{FATHER}[v_i] = 0$ 。初始时, 每个顶点 v_i 本身是一棵树。为了定位集合 i 的根, 构建数组 ROOT , 其中 $\text{ROOT}[i]$ 指向 v_i 。为定位对应元素 i 的顶点, 构建数组 ELEMENT , 初始时, $\text{ELEMENT}[i] = v_i$ 。

2. **执行 FIND(i)。** 图 4-18 给出相应的程序。由顶点 $\text{ELEMENT}[i]$ 开始, 沿着指向树根的路径, 将碰到的顶点置于一个列表中。到达根时, 打印集合的名字, 且使在遍历路径上的每个顶点都成为根的儿子。

3. 执行 $\text{UNION}(i, j, k)$ 。通过数组 ROOT ，找到表示集合 i 和 j 的树根。然后，使较小树的根成为较大树的根的儿子。如图 4-19。 □

```

begin
  make LIST empty;
  v ← ELEMENT[i];
  while FATHER[v] ≠ 0 do
    begin
      add v to LIST;
      v ← FATHER[v]
    end;
  comment v is now the root;
  print NAME[v];
  for each w on LIST do FATHER[w] ← v
end

```

图 4-18 执行指令 $\text{FIND}(i)$

```

begin
  wlg assume COUNT[ROOT[i]] ≤ COUNT[ROOT[j]]
  otherwise interchange i and j in
  begin
    LARGE ← ROOT[j];
    SMALL ← ROOT[i];
    FATHER[SMALL] ← LARGE;
    COUNT[LARGE] ← COUNT[LARGE] + COUNT[SMALL];
    NAME[LARGE] ← k;
    ROOT[k] ← LARGE
  end
end

```

图 4-19 执行指令 $\text{UNION}(i, j, k)$

接下来证明，路径压缩法相当大地加快了算法的执行。为了计算这一改进，引入函数 F 和 G ，使得

$$F(0) = 1$$

$$F(i) = 2^{F(i-1)}, \text{ 对于 } i > 0.$$

如图 4-20 中的表所示，函数 F 的增长非常快。函数 $G(n)$ 定义为满足 $F(k) \geq n$ 的最小整数 k 。函数 G 的增长则非常慢。实际上，对于所有“实际的” n 值，即对所有的 $n \leq 2^{65536}$ ， $G(n) \leq 5$ 。

现在证明算法 4.3 可以在至多 $c'nG(n)$ 时间内，执行由 cn 条 UNION 和 FIND 指令组成的指令序列 σ ，其中， c 和 c' 为常量， c' 依赖于 c 。为了简单，假定执行一条 UNION 指令消耗一个“时间单元”，执行指令 $\text{FIND}(i)$ 消耗的时间单元数从标记为 i 的顶点到包含该顶点的树根的路径上的顶点数成比例[⊖]。

定义 对于 UNION 和 FIND 指令组成的指令序列 σ ，可以如下定义一个顶点的秩(rank)：

n	$F(n)$
0	1
1	2
2	4
3	16
4	65536
5	2^{65536}

图 4-20 F 中的值

⊖ 在此，“时间单元”对应 RAM 所需的常数步。由于忽略了常量因子，结果数量级也可以用“时间单元”来表示。

1. 从 σ 中删除 FIND 指令;
2. 执行由 UNION 指令组成的结果序列 σ' ;
3. 在最后得到的森林中, 顶点 v 的秩就是 v 的高度。

现在来导出一些与顶点的秩有关的重要性质。

引理 4.2 存在至多 $n/2^r$ 个秩为 r 的顶点。

证明: 由引理 4.1, 在执行 σ' 得到的森林中, 秩为 r 的每个顶点至少有 2^r 个后代。由于森林中任何两个相同高度的不同顶点后代的集合是不相交的, 并且 2^r 个或者更多的顶点至多存在 $n/2^r$ 个不相交的集合, 所以至多有 $n/2^r$ 个秩为 r 的顶点。 \square

推论 不存在秩超过 $\log n$ 的顶点。

引理 4.3 在执行 σ 的某个时刻, 如果 w 是 v 的某个后代, 则 w 的秩小于 v 的秩。

证明: 如果在执行 σ 的某个时刻, w 成为 v 的后代, 那么在执行 σ' 得到的森林中, w 将是 v 的后代。因此, w 的高度必定小于 v 的高度, 因此 w 的秩小于 v 的秩。 \square

现在将秩划分为组(group)。将秩 r 放入组 $G(r)$ 中。例如, 将秩 0 和 1 放入组 0, 秩 2 放入组 1, 秩 3 和 4 放入组 2, 秩 5 ~ 16 放入组 3。对于 $n > 1$, 最大可能的秩 $\lfloor \log n \rfloor$ 位于秩组(rank group) $G(\lfloor \log n \rfloor) \leq G(n) - 1$ 中。

考虑执行由 cn 条 UNION 和 FIND 指令组成的指令序列 σ 的代价。执行每条 UNION 指令为一个时间单位的代价, σ 中的所有 UNION 指令可在 $O(n)$ 时间内执行。为了限定执行所有 FIND 指令的代价, 使用一个重要的“簿记”(bookkeeping)技术。执行单条 FIND 指令的代价在 FIND 指令本身和在森林数据结构中路径上的实际移动过的顶点之间进行分摊。通过对分摊给所有 FIND 指令的代价进行求和, 并加上分摊给森林中所有顶点的代价, 可以计算出总的代价。

如下计量执行 FIND(i) 指令的代价。设 v 为从表示 i 的顶点到包含 i 的树根的路径上的一个顶点。

1. 如果 v 为根, 或者 FATHER[v] 处于一个不同于 v 的秩组中, 则为 FIND 指令本身计量一个时间单位;
2. 如果 v 及其父亲处于相同的秩组中, 则计量 v 承担一个时间单位。

由引理 4.3, 沿着路径往上, 顶点的秩单调增加, 由于存在至多 $G(n)$ 个不同的秩组, 所以在规则 1 下, 不会有 FIND 指令承担超过 $G(n)$ 个时间单位。如果应用规则 2, 顶点 v 将被移动, 成为一个秩大于其父的顶点的儿子。如果顶点 v 在秩组 $g > 0$ 中, 则在 v 得到位于更高秩组中的父亲之前, 可以移动 v , 并计量至多 $F(g) - F(g-1)$ 次。在秩组 0 中, 在得到位于更高秩组中的父亲之前, 可移动一个顶点至多一次。由此, 根据规则 1, 移动 v 的代价将由 FIND 指令来承担。

为了得到顶点本身的代价的上界, 可将秩组中顶点的最大可能代价乘以秩组中的顶点数, 并对所有秩组求和。设 $N(g)$ 是秩组 $g > 0$ 中的顶点数, 由引理 4.2:

$$\begin{aligned}
 N(g) &\leq \sum_{r=F(g-1)+1}^{F(g)} n/2^r \\
 &\leq (n/2^{F(g-1)+1}) \left[1 + \frac{1}{2} + \frac{1}{4} + \cdots \right] \\
 &\leq n/2^{F(g-1)} \\
 &\leq n/F(g)
 \end{aligned}$$

在秩组 $g > 0$ 中, 任意顶点的最大代价小于等于 $F(g) - F(g-1)$ 。因此, 在秩组 g 中, 所有顶点的最大代价界为 n 。很明显, 同样的讨论也适用于 $g = 0$ 的情形。由于存在至多 $G(n)$ 个秩组, 所有顶点的最大代价为 $nG(n)$ 。因此, 处理 cn 条 FIND 指令需要的时间总数至多 $cnG(n)$ 次分摊到

FIND 指令, 至多 $nG(n)$ 次分摊给顶点。因此, 有下列定理。

定理 4.4 设 c 为任意常数, 则存在另一个依赖于 c 的常数 c' , 使算法 4.3 能在至多 $c'nG(n)$ 个时间单位内, 执行由 cn 条对 n 个元素的 UNION 和 FIND 指令组成的指令序列 σ 。

证明: 由上述讨论可证。□

如果在序列 σ 中, 允许有基本操作 INSERT、DELETE 以及 UNION 和 FIND 指令, 则 σ 仍然能够在 $O(nG(n))$ 时间内执行。证明留作练习。

目前还不知道定理 4.4 是否给算法 4.3 的运行时间提供了一个紧界(tight bound)。然而, 出于理论上的重要性, 本节的余下部分将证明, 算法 4.3 的运行时间和 n 不成线性关系。为此, 构建一个特定的 UNION 和 FIND 指令序列, 算法 4.3 需要多于线性的时间来处理该序列。

为了方便, 引入一个基于树的新操作, 称为局部 FIND(partial FIND), 或 PF。设 T 为一棵树, 其中 $v, v_1, v_2, \dots, v_m, w$ 是从顶点 v 到祖先 w 的一条路径(w 不一定为根)。操作 $PF(v, w)$ 使 $v, v_1, v_2, \dots, v_{m-1}$ 成为顶点 w 的儿子。称该局部 FIND 有长度 $m+1$ (若 $v=w$, 则长度为 0)。图 4-21b 给出了对于图 4-21a 中的树执行 $PF(v, w)$ 后的效果。

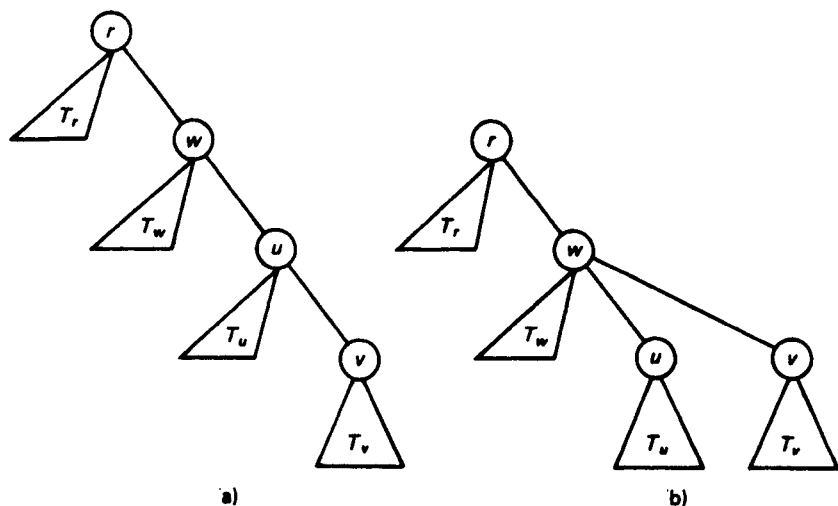


图 4-21 局部 FIND 操作的效果

假设给定一个 UNION 和 FIND 指令序列 σ , 当在 σ 中执行一条给定的 FIND 指令时, 在某棵树 T 上定位顶点 v , 并跟踪从 v 到树 T 的根 w 的路径。现在假定仅执行 σ 中的 UNION 指令, 而不考虑 FIND 指令, 其结果为树的森林 F 。通过在 F 中定位由最初的 FIND 指令所使用的顶点 v 和 w , 然后执行 $PF(v, w)$, 仍然可以得到执行 σ 中一条给定 FIND 指令的效果。注意, 顶点 w 可能不再是 F 中的根。

为获取算法 4.3 运行时间的下界, 基于 UNION 指令序列来考虑算法的性能, 接着执行 PF, 这些 PF 操作可以被执行时间相同的 UNION 和 FIND 指令所组成的指令序列所取代。从下面特定的树中, 能够得到特定的 UNION 和 PF 序列, 其中, 算法 4.3 所消耗的时间多于线性时间。

定义 对于 $k \geq 0$, 设 $T(k)$ 为树, 满足

1. $T(k)$ 的每个叶子的深度为 k ;
2. 高为 h 的每个顶点都有 2^h 个儿子, $h \geq 1$ 。

因此, $T(k)$ 的根有 2^k 个儿子, 每个儿子都是 $T(k-1)$ 副本的根。图 4-22 给出树 $T(2)$ 。

引理 4.4 对于任意的 $k \geq 0$, 通过 UNION 指令序列, 可以创建一棵树 $T'(k)$, 包含在树

$T(k)$ 的一个子图中。并且,在 $T'(k)$ 中,至少有四分之一的顶点是 $T(k)$ 中的叶子。

证明:对 k 进行归纳。对于 $k=0$,因为 $T(0)$ 仅包含一个顶点,引理显然正确。对于 $k>0$,要构造 $T'(k)$,首先构建 2^k+1 个 $T'(k-1)$ 的副本。选择 $T'(k-1)$ 的一个副本组成树 $T'(k)$,然后将其他的副本逐个合并进来。最终得到的树根有 2^k 个儿子,每个都是 $T'(k-1)$ 的根。

设 $N'(k)$ 为 $T'(k)$ 中的顶点总数, $L(k)$ 是 $T(k)$ 中的叶子数。则

$$N'(0) = 1$$

$$N'(k) = (2^k + 1)N'(k-1), \text{ 对于 } k \geq 1$$

并且

$$L(0) = 1$$

$$L(k) = 2^k L(k-1), \text{ 对于 } k \geq 1$$

因此

$$\frac{L(k)}{N'(k)} = \frac{\prod_{i=1}^k 2^i}{\prod_{i=1}^k (2^i + 1)} = \frac{2}{3} \prod_{i=2}^k \frac{1}{1 + 2^{-i}} \quad \text{对于 } k \geq 1 \quad (4-3)$$

注意,对于 $i \geq 2$,有 $\ln(1 + 2^{-i}) < 2^{-i}$,因此

$$\ln\left(\prod_{i=2}^k \frac{1}{1 + 2^{-i}}\right) \geq -\sum_{i=2}^k 2^{-i} \geq -\frac{1}{2} \quad (4-4)$$

通过式(4-3)和式(4-4),有

$$\frac{L(k)}{N'(k)} \geq \frac{2}{3} e^{-1/2} \geq \frac{1}{4}$$

因此,定理得证。□

考虑构造 UNION 和 PF 指令序列。首先构建树 $T'(k)$,接着在子图 $T(k)$ 的叶子上执行 PF 操作。现在证明,对于每个 $l>0$,存在一个 k ,能够在 $T(k)$ 的每个叶子上连续执行长度为 l 的 PF 操作。

定义 设 $D(c, l, h)$ 为 k 的最小值,如果用任意有 l 个叶子且高度至少为1的树取代根的高度为 h 的 $T(k)$ 中的每一棵子树,那么可以对所得树的每个叶子执行长度为 c 的 PF 操作。

引理 4.5 对于所有大于0的 c, l 和 h ,定义 $D(c, l, h)$ (即它是有限的)。

证明:包含两层归纳。我们希望通过 c 的归纳来证明结果。但是,给定 $c-1$ 的结果,要证明 c 的结果,还必须对 l 进行归纳。

归纳基础 $c=1$ 是容易的。对长度为1的 PF 操作,它不会移动任何顶点,因此,对于所有的 l 和 h , $D(1, l, h) = h$ 。

现在对 c 进行归纳。假定对所有的 l 和 h ,定义 $D(c-1, l, h)$ 。必须证明对所有的 l 和 h ,可定义 $D(c, l, h)$ 。为此,对 l 进行归纳。

对归纳基础,证明

$$D(c, l, h) \leq D(c-1, 2^{h+1}, h+1)$$

注意,当 $l=1$ 时,用子树的一片叶子取代树,这些子树的根是 $T(k)$ (对某些 k)中高度为 h 的顶点。设 H 表示该 $T(k)$ 中高度为 h 的顶点集合。显然,在修改后的树中,每片叶子都是 H 中某一成员的原有后代。因此,如果能够对 H 的所有成员执行长度为 $c-1$ 的 PF 操作,则肯定可以对所有叶子执行长度为 c 的 PF 操作。

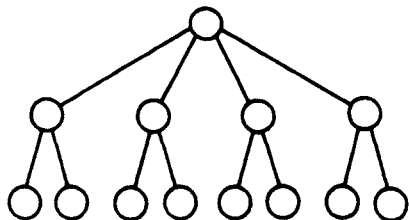


图 4-22 树 $T(2)$

设 $k = D(c - 1, 2^{h+1}, h + 1)$ 。由基于 c 的归纳假设, 可知 k 存在。如果考虑 $T(k)$ 中高度为 $h + 1$ 的顶点, 可以看到每个顶点都有 2^{h+1} 个儿子, 且都是 H 的成员。若在 $T(k)$ 中删除 H 中顶点的所有原有后代, 实际上是对每一棵根的高度为 $h + 1$ 的子树, 用 2^{h+1} 片叶子取代高度为 1 的树。由 D 的定义, $k = D(c - 1, 2^{h+1}, h + 1)$ 足够大, 以致可以对所有叶子 (即 H 的成员) 执行长度为 $c - 1$ 的 PF 操作。

现在, 为完成对 c 的归纳, 必须对 l 进行归纳。特别地, 我们将证明:

$$D(c, l, h) \leq D(c - 1, 2^{D(c, l-1, h)(1 + D(c, l-1, h))/2}, D(c, l-1, h)), \text{ 对于 } l > 1 \quad (4-5)$$

要证明式 (4-5), 设 $k = D(c, l - 1, h)$, k' 为式 (4-5) 的右边。对于 $T(k')$ 中每个高度为 h 的顶点, 必须找到一种方式替换一棵包含 l 片叶子的树, 然后对每片叶子执行长度为 c 的 PF 操作。首先, 对要取代树中的叶子, 取 $l - 1$ 来执行 PF 操作。由对 l 归纳的归纳假设, 可以对这些子树中被取代树的叶子, 取 $l - 1$ 来执行 PF 操作。

在对叶子的 $l - 1$ 执行 PF 操作后, 发现每棵被取代树的第 l 片叶子现在都有一个不同于任何其他被取代树的第 l 片叶子的父亲。称这类父亲的集合为 F 。如果可以对父亲执行长度为 $c - 1$ 的 PF 操作, 则可以对叶子执行长度为 c 的 PF 操作。设 S 为 $T(k')$ 中的一棵子树, 其根的高度为 k 。很容易验证, 在 $T(k')$ 中, S 有 $2^{k(k+1)/2}$ 片叶子。因此, 在执行了 PF 操作之后, 在 S 中同时也在 F 中的顶点数至多为 $2^{k(k+1)/2}$ 。因此, S 的余下部分可认为是有 $2^{k(k+1)/2}$ 叶子的一棵任意树, 顶点在 F 中。由 c 和 l 的归纳假设, 式 (4-5) 得证。□

定理 4.5 算法 4.3 的时间复杂度大于 cn , 其中 c 是常数。

证明: 假定存在常数 c , 使算法 4.3 可以在不超过 cn 个时间单位内, 执行由 $n - 1$ 条 MERGE 和 n 条 FIND 指令组成的任意指令序列。选择 $d > 4c$, 计算 $k = D(d, 1, 1)$ 。通过 UNION 指令序列, 构造 $T'(k)$ 。由于可以对要嵌入的树 $T(k)$ 的每片叶子执行长度为 d 的 PF 操作, 且 $T(k)$ 的叶子由超过 $1/4$ 的 $T'(k)$ 中的顶点组成, 该 UNION 和 PF 指令序列将需要多于 cn 个时间单位, 矛盾。□

4.8 UNION-FIND 算法的应用和扩展

现在已经知道, 如何很自然地由例 4-1 的生成树问题中得出由基本指令 UNION 和 FIND 所组成的操作序列。本节给出其他一些使用 UNION 和 FIND 指令序列的问题。第一个问题涉及离线计算, 也就是说, 在得到任何答案之前, 可以读取整条指令序列。

应用 1 离线 MIN 问题

给定指令 INSERT(i) 和 EXTRACT_MIN, 从初始为空的集合 S 开始。每当碰到 INSERT(i) 指令, 则将整数 i 放入集合 S 中。每次执行 EXTRACT_MIN 指令时, 在 S 中查找最小元素, 并删除。

设 σ 为包含 INSERT 和 EXTRACT_MIN 的指令序列, 满足: 对每个 $i (1 \leq i \leq n)$, INSERT(i) 指令最多出现一次。给定序列 σ , 要找出被 EXTRACT_MIN 指令删除的整数系列。甚至在需要计算输出序列的第一个元素之前, 就已经给定了完整的序列 σ , 所以该问题是离线的。

可通过下述方法来解决离线 MIN 问题。设 k 为 σ 中 EXTRACT_MIN 指令的条数。可以将 σ 写为 $\sigma_1 E \sigma_2 E \sigma_3 E \cdots \sigma_k E \sigma_{k+1}$, 其中, 对于 $1 \leq j \leq k + 1$, 每个 σ_j 仅包含 INSERT 指令, E 代表一条 EXTRACT_MIN 指令。下面通过算法 4.3 的集合合并算法模拟 σ 。对此集合合并算法, 如果 INSERT(i) 指令出现在子序列 σ_j 中, 则通过让名字为 $j (1 \leq j \leq k + 1)$ 的集合包含元素 i 来初始化一个集合序列。若存在名字为 j 的集合, 则通过使用两个数组 PRED 和 SUCC, 为 j 中的值创建一个双链接的有序表。初始时, 对于 $1 \leq j \leq k + 1$, PRED[j] = $j - 1$, 对于 $0 \leq j \leq k$, SUCC[j] = $j + 1$ 。然

后, 执行图 4-23 的程序。

```

for i ← 1 until n do
  begin
    j ← FIND(i);
    if j ≤ k then
      begin
        print i "is deleted by the "j"th EXTRACT_MIN instruction";
        UNION(j, SUCC[j], SUCC[j]);
        SUCC[PRED[j]] ← SUCC[j];
        PRED[SUCC[j]] ← PRED[j]
      end
    end
  end
end

```

图 4-23 离线 MIN 问题的程序

显然, 该程序的执行时间受制于集合合并算法的运行时间。因此, 离线 MIN 问题的时间复杂度为 $O(nG(n))$ 。

例 4.8 考虑指令序列 $\sigma = 43E2E1E$, 其中 j 代表 $\text{INSERT}(j)$, E 代表 EXTRACT_MIN 。因此, $\sigma_1 = 43$, $\sigma_2 = 2$, $\sigma_3 = 1$, σ_4 为空序列。初始的数据结构为集合序列

$$1 = \{3, 4\} \quad 2 = \{2\} \quad 3 = \{1\} \quad 4 = \emptyset$$

首次执行 **for** 循环, 确定 $\text{FIND}(1) = 3$ 。因此, σ 中第 3 条 EXTRACT_MIN 指令的结果为 1。集合序列变为

$$1 = \{3, 4\} \quad 2 = \{2\} \quad 4 = \{1\}$$

此时, 由于名字为 3 和 4 的集合合并成了名字为 4 的单个集合, $\text{SUCC}[2]$ 对应等于 4 的集合, $\text{PRED}[4]$ 对应等于 2 的集合。

接着用 $i=2$ 执行下一步, 确定 $\text{FIND}(2) = 2$ 。因此, 第 2 条 EXTRACT_MIN 指令的结果为 2。将名字为 2 的集合同其后继集合(名字为 4)合并, 得到集合序列

$$1 = \{3, 4\} \quad 4 = \{1, 2\}$$

最后两步确定, 第 1 条 EXTRACT_MIN 指令的结果为 3, 而 4 则始终没有得到抽取。□

下个应用是深度测定问题, 它常出现在汇编语言程序中标识符“等价性”的问题中。许多汇编语言允许声明语句中的两个标识符表示相同存储位置。如果汇编器碰到一条语句声明两个等价标识符 α 和 β , 那么它必须找到两个集合 S_α 和 S_β , 分别表示等价于 α 和 β 的标识符集, 并用它们的并取代这两个集合。显然, 该问题可用 UNION 和 FIND 指令序列来模拟。

然而, 如果更仔细地研究该问题, 可以找到另一种方法来应用前一节的数据结构。在符号表中, 每个标识符都有一项, 如果一组标识符等价, 则可以方便地将与它们有关的数据仅保存在其中一个符号表项处。这表明, 对于每个等价的标识符集, 存在一个源点(origin), 用来保存符号表中与集合有关的信息, 每个成员都有一个与源点的偏移。将标识符的偏移加到集合的源点, 可以找到该标识符在符号表中的对应位置。然而, 要使两个标识符集等价, 还必须修改与源点的偏移量。应用 2 概要地描述了偏移的修改问题。

应用 2 深度测定问题

给定包含两类指令 $\text{LINK}(v, r)$ 和 $\text{FIND_DEPTH}(v)$ 的指令序列。初始时, 有 n 棵无向有根树, 每棵包含唯一的顶点 $i (1 \leq i \leq n)$ 。指令 $\text{LINK}(v, r)$ 其中 r 为一棵树的根, 且 v 为不同的树上的顶点, 使得根 r 成为顶点 v 的一个儿子。 v 和 r 位于不同树中且 r 为根, 这些条件确保最终得到的图仍然是一个森林。指令 $\text{FIND_DEPTH}(v)$ 需要确定并打印顶点 v 的当前深度。

如果用一个常规的邻接表表示维护该森林,并以显见的方法测定顶点的深度,则该算法的复杂度将是 $O(n^2)$ 。但是,这里将考虑使用称为 D 森林的结构代替原先的森林结构。 D 森林的唯一目的是为了可以快速计算深度。给 D 森林中的每个顶点设置一个整数权重,使得沿 D 森林中从顶点 v 到根的路径的总权重为 v 在初始森林中的深度。对于 D 森林中的每棵树,计算树中的顶点数目。

起初, D 森林包含 n 棵树,每棵树包含与唯一的整数 $i (1 \leq i \leq n)$ 对应的单一顶点。每个顶点的初始权重为 0。

沿着从 v 到根 r 的路径执行指令 $\text{FIND_DEPTH}(v)$ 。设 v_1, v_2, \dots, v_k 是路径 ($v_1 = v, v_k = r$) 上的顶点,则

$$\text{DEPTH}(v) = \sum_{i=1}^k \text{WEIGHT}[v_i]$$

接着进行路径压缩。使每个 $v_i (1 \leq i \leq k-2)$ 成为根 r 的儿子。为维护权重属性,对于 $1 \leq i < k$, 顶点 v_i 的新权重 (WEIGHT) 必须变成 $\sum_{j=i}^{k-1} \text{WEIGHT}[v_j]$ 。由于可在 $O(k)$ 时间内计算新权重,一条 FIND_DEPTH 指令的执行时间与 FIND 指令有相同的数量级。

指令 $\text{LINK}(v, r)$ 的执行结果是合并含有顶点 v 和 r 的树,再将较小的树合并到更大的树。设 T_v 和 T_r 分别为 D 森林中包含 v 和 r 的树,且 v' 和 r' 分别为 T_v 和 T_r 的根。该 D 森林中的树没有必要同初始森林中的树同构,在特殊情形下, r 可能不是 T_r 的根。用 $\text{COUNT}(T)$ 表示树 T 中的顶点数。需要考虑两种情形。

情形 1 $\text{COUNT}(T_v) \leq \text{COUNT}(T_r)$ 。使 r' 为 v' 的儿子。必须调整 T_r 原先的根 r' 的权重,以便沿着合并后的树中从 w 到 v' 的新路径正确计算 T_r 中每个顶点 w 的深度。为此,执行指令 $\text{FIND_DEPTH}(v)$, 然后如下操作:

$$\text{WEIGHT}[r'] \leftarrow \text{WEIGHT}[r'] - \text{WEIGHT}[v'] + \text{DEPTH}(v) + 1$$

因此,通过使 v 的深度加 1,使 T_r 中每个顶点的深度有效增加。最终,使合并树的计数等于 T_v 和 T_r 的计数之和;

情形 2 $\text{COUNT}(T_v) > \text{COUNT}(T_r)$ 。执行 $\text{DEPTH}(v)$, 而后使 v' 成为 r' 的儿子,并如下操作:

$$\text{WEIGHT}[r'] \leftarrow \text{WEIGHT}[r'] + \text{DEPTH}(v) + 1$$

$$\text{WEIGHT}[v'] \leftarrow \text{WEIGHT}[v'] - \text{WEIGHT}[r']$$

$$\text{COUNT}(T_v) \leftarrow \text{COUNT}(T_v) + \text{COUNT}(T_r)$$

总之,可以在 $O(nG(n))$ 时间内执行 $O(n)$ 条 LINK 和 FIND_DEPTH 指令。

应用 3 有穷自动机的等价性

确定性有穷自动机是能够识别符号串的一种机器。如图 4-24 所示,它包含一条以方框表示的输入带、一个输入磁头以及一个有穷状态控制器。用 5 元组 (S, I, δ, s_0, F) 表示有穷自动机 M , 其中:

1. S 为有穷控制器 (finite control) 的状态集;
2. I 为输入字母表 (input alphabet)。输入带上的每个方框都包含 I 中的一个符号;
3. δ 是从 $S \times I$ 到 S 的映射。如果 $\delta(s, a) = s'$, 则 M 处于状态 s 且输入头处在符号 a 处时, M 将输入头右移,并进入状态 s' ;
4. s_0 是 S 中的一个特定状态,称为开始状态 (start state);
5. F 是 S 中的一个特定子集,称为接收 (或终止) (accepting or fi-

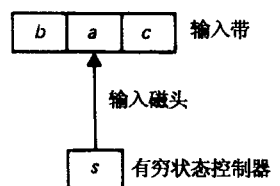


图 4-24 有穷自动机

nal) 状态集。

用 I^* 表示 I 中所有有限长的符号串的集合。 I^* 包含空串 ϵ 。按以下方式扩展函数 δ , 将 $S \times I^*$ 映射到 S :

1. $\delta(s, \epsilon) = s$;
2. 对 I^* 中的所有 x 和 I 中的 a , $\delta(s, xa) = \delta(\delta(s, x), a)$ 。

如果 $\delta(s_0, x) \in F$, 则输入串 x 为 M 所接收。 M 接收的语言记为 $L(M)$, 其为 M 接收的字符串集合。9.1 节包含更广泛的对无穷自动机的介绍。

如果对于 I^* 中的每个 x , $\delta(s_1, x)$ 是一个接收状态当且仅当 $\delta(s_2, x)$ 是一个接收状态, 则称状态 s_1 和 s_2 是等价的。

如果 $L(M_1) = L(M_2)$, 则称两个有穷自动机 M_1 和 M_2 是等价的。本节将证明可以使用 UNION-FIND 算法在 $O(nG(n))$ 步内测定两个有穷自动机 $M_1 = (S_1, I, \delta_1, s_1, F_1)$ 和 $M_2 = (S_2, I, \delta_2, s_2, F_2)$ 是否等价, 其中, $n = \|S_1\| + \|S_2\|$ 。

状态等价的一个重要属性在于, 如果两个状态 s 和 s' 等价, 那么对于所有的输入符号 a , $\delta(s, a)$ 和 $\delta(s', a)$ 也必须等价。由于存在空串, 接收态不可能等价于非接收态。因此, 如果初始假定 M_1 和 M_2 的开始状态 s_1 和 s_2 等价, 则可以推断某些其他的状态对必定也等价。如果这些状态对中的一个包含接收态和非接收态, 则 s_1 和 s_2 不等价。尽管没有给出证明, 但是反之亦然。

为测定两个有穷自动机 $M_1 = (S_1, I, \delta_1, s_1, F_1)$ 和 $M_2 = (S_2, I, \delta_2, s_2, F_2)$ 是否等价, 可按如下方式来处理:

1. 假定两个初始状态 s_1 和 s_2 等价, 可以用图 4-25 的算法确定必定等价的所有状态集。LIST 保存状态对 (s, s') , 其中 s 和 s' 已被发现是等价的, 但是其后继 $(\delta(s, a), \delta(s', a))$ 的等价性尚未测定。初始, LIST 仅包含初始状态对 (s_1, s_2) 。为找出等价状态集, 程序使用上述关于不相交集合并算法。COLLECTION 表示集合族。初始时, $S_1 \cup S_2$ 中的每个状态都处于其自身的集合中 (不失一般性, 可以假定 S_1 和 S_2 中的状态是不相交的)。当发现两个状态 s 和 s' 等价, 合并 COLLECTION 中包含 s 和 s' 的集合 A 和 A' , 并将得到的集合称为 A 。

2. 程序的其余部分将 COLLECTION 中的集合表示为一个划分, 该划分将 $S_1 \cup S_2$ 划分为必定等价的状态块。 M_1 和 M_2 是等价的, 当且仅当没有既包含接收态也包含非接收态的块。

```

begin
  LIST  $\leftarrow (s_1, s_2)$ ;
  COLLECTION  $\leftarrow \emptyset$ ;
  for each  $s$  in  $S_1 \cup S_2$  do add  $\{s\}$  to COLLECTION;
  comment 对  $S_1 \cup S_2$  中的每个状态初始化一个集合;
  while there is a pair  $(s, s')$  of states on LIST do
    begin
      delete  $(s, s')$  from LIST;
      let  $A$  and  $A'$  be FIND( $s$ ) and FIND( $s'$ ), respectively;
      if  $A \neq A'$  then
        begin
          UNION( $A, A', A$ );
          for all  $a$  in  $I$  do
            add  $(\delta(s, a), \delta(s', a))$  to LIST
          end
        end
      end
    end
end

```

图 4-25 假定 s_1 和 s_2 等价时, 找出等价状态集合的算法

集合合并算法占据了大部分算法的运行时间(作为状态数 $n = \|S_1\| + \|S_2\|$ 的函数)。由于每条 UNION 指令都使 COLLECTION 中的集合数减少 1, 可能存在至多 $n-1$ 个 UNION, 并且开始时仅有 n 个集合。FIND 的数目同 LIST 中状态对的数目成比例。由于除了初始状态对 (s_1, s_2) , 仅在一条 UNION 指令之后将一对状态放置到 LIST 中, 因此, FIND 的数目至多为 $n \times \|I\|$ 。由此, 假定 $\|I\|$ 为常量, 则确定 M_1 是否等价于 M_2 所需要的时间为 $O(nG(n))$ 。

4.9 平衡树方案

有几类重要的问题, 虽然它们同 UNION-FIND 问题相似, 但是讨论这些问题仍需要回到前面用 $O(n \log n)$ 时间(最坏情形)处理 n 条指令序列的情况。这类问题之一是处理由 MEMBER、INSERT 和 DELETE 指令所组成的指令序列, 其中总的可能元素比实际所用元素的数量更大。在这种情况下, 无法通过直接索引到指针数组来访问一个元素。必须用散列表, 或者用二叉查找树。

如果插入 n 个元素, 散列方法所需的平均访问时间为常量, 但是最坏情况下每次访问的时间为 $O(n)$ 。使用二叉查找树时, 每次访问的预期访问时间为 $O(\log n)$ 。但是, 倘若名字集合不是静态的, 二叉查找树也可能会得到一个糟糕的最坏情况访问时间。如果只是简单地将名字增加到树中, 没有机制来保持其平衡, 则最终可能得到一棵包含 n 个元素但深度接近于 n 的树。因此, 在最坏情况下二叉查找树每个操作的执行时间可能是 $O(n)$ 。本节提供的技术将使最坏情况的执行时间减少为每个操作 $O(\log n)$ 步。

另一类需要 $O(n \log n)$ 时间的问题用于在线处理 n 条指令的序列, 其中包含 INSERT、DELETE 和 MIN 操作。如果需要表示有序列表并且执行列表的链接和分离, 则将产生第 3 类类似的问题。

本节的技术将允许在线处理 4.1 节所提及的集合上 7 类基本操作的特定子集所组成的指令序列。所用基本数据结构为平衡树(balanced tree), 树的高度约等于结点总数的对数值。

初始平衡树是很容易构造的。然而, 问题在于当执行 INSERT 和 DELETE 指令序列时, 如何防止树变成不平衡的。例如, 若执行 DELETE 命令序列, 仅从树的左边删除顶点, 最终得到的树将左右不平衡。解决方法是周期性地重新平衡这棵树。

许多机制可用在必要的时候重新平衡树。其中的几种方法是让树结构足够灵活, 使高为 h 的树中顶点数处于 2^h 到 2^{h+1} 或 3^h 的范围内。在必须对子树的根做任何改变之前, 这些方法允许子树中的顶点数至少增加为原来的两倍。

接下来讨论具备这种性质的两类特定方法: 2-3 树和 AVL 树。处理 2-3 树的算法在概念上更易于理解, 后面将进行讨论。处理 AVL 树的算法与 2-3 树相似, 可在习题中找到。

定义 2-3 树(2-3tree)是一棵树, 每个非叶子的结点有 2 或 3 个儿子, 且从根到叶子的每条路径有相同的长度。注意, 只包含一个结点的树是 2-3 树。图 4-26 给出两棵包含 6 片叶子的 2-3 树。

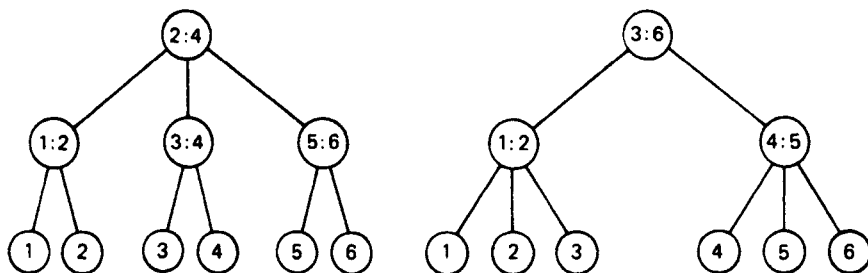


图 4-26 2-3 树

下列引理给出了 2-3 树中顶点和叶子数同其高度之间的关系。

引理 4.6 设 T 为一棵高 h 的 2-3 树。 T 中的顶点数位于 $2^{h+1} - 1$ 和 $(3^{h+1} - 1)/2$ 之间, 并且叶子数位于 2^h 和 3^h 之间。

证明: 对 h 进行基本归纳。 □

通过把集合中的元素指派给树的叶子可以用一棵 2-3 树表示一个线性有序集 S 。使用 $E[l]$ 表示存储在叶子 l 中的元素。存在两种基本的将元素指派给叶子的方法, 采用哪种方法依赖于具体应用。

如果可能的元素数目远远大于所用到的实际元素数量, 且将树用作为字典, 那么可以从左到右按升序指派元素。对每个非叶顶点 v , 还需两类其他信息 $L[v]$ 和 $M[v]$ 。 $L[v]$ 是已分配给子树的 S 中的最大元素, 该子树的根为 v 的最左儿子; $M[v]$ 为已分配给子树的 S 中的最大元素, 该子树的根为 v 的第 2 个儿子。参看图 4-26。与顶点相关联的 L 和 M 值使得可以从根开始, 以一种类似于二分搜索的方式查找一个元素。找到任何元素的时间同树的高度成比例。因此, 如果使用具有该性质的一棵 2-3 树表示集合, 则可以在 $O(\log n)$ 时间内处理拥有 n 个元素的集合上的一条 MEMBER 指令。

将元素指派到叶子的第 2 种方法是不对元素分配的顺序设限。在实现指令 UNION 的时候, 该方法特别有用。然而, 为执行诸如 DELETE 这样的指令, 则需要有一种机制来定位表示给定元素的叶子。如果集合中的元素是某一确定范围内的整数, 假定 $1 \sim n$, 则通过某个数组的第 i 个位置可以找到表示元素 i 的叶子。如果集合中的元素来自于某个大的全集, 那么可以通过使用备用字典来找到表示元素 i 的叶子。

考虑下面指令集:

1. INSERT, DELETE, MEMBER
2. INSERT, DELETE, MIN
3. INSERT, DELETE, UNION, MIN
4. INSERT, DELETE, FIND, CONCATENATE, SPLIT

一般称用来处理集合 1 中指令的数据结构为字典; 称可处理集合 2 中指令的数据结构为优先队列 (priority queue); 称处理集合 3 中指令的数据结构为可合并堆 (mergeable heap); 称处理集合 4 中指令的数据结构为可连接队列 (concatenable queue)。

下面证明 2-3 树可用来实现字典、优先队列、可连接队列以及可合并堆, 并可在 $O(n \log n)$ 时间内处理 n 条指令。所用的技术也完全能够执行由本章开始时所列出的 7 条指令的任何可兼容子集组成的指令序列。唯一的不兼容性体现在: UNION 意味着无序集, 而 SPLIT 和 CONCATENATE 则意味着有序集。

4.10 字典和优先队列

本节将考虑实现字典和优先队列所需要的基本操作。假定按从左到右的顺序将元素指派给 2-3 树中的叶子, 且对于每个非叶顶点 v , 有 $L[v]$ 和 $M[v]$ (前一节所描述的“最大”后继函数)。

为插入新元素 a 到一棵 2-3 树中, 必须确定 a 所对应的新叶子 l 的位置。为此, 尝试在树中定位元素 a 。假定树包含不止一个元素, 则对 a 的查找终止于顶点 f , 其中 f 包含两或三片作为叶子的叶子。

如果 f 仅有两片叶子 l_1 和 l_2 , 则使 l 成为 f 的儿子。如果 $a < E[l_1]^\ominus$, 则使 l 成为 f 的最左儿子, 并设定 $L[f] = a$ 且 $M[f] = E[l_1]$; 如果 $E[l_1] < a < E[l_2]$, 则使 l 成为 f 的中间儿子, 并设定

$\ominus E[v]$ 为存储在叶子 v 上的元素。

$M[f] = a$; 如果 $E[l_2] < a$, 则使 l 成为 f 的第 3 个儿子。在后一种情形, 可能需要改变 f 某些祖先的 L 或 M 值。

例 4.9 如果插入元素 2 到图 4-27a 所示的 2-3 树中, 则得到图 4-27b 所示的 2-3 树。 □

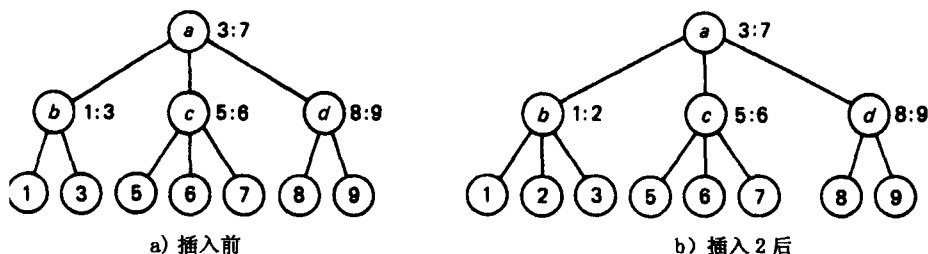


图 4-27 插入元素到一棵 2-3 树

现在假定 f 已经有 3 片叶子 l_1 、 l_2 和 l_3 。使 l 成为 f 中合适位置的儿子。现在顶点 f 有 4 个儿子。为了保持 2-3 树的性质, 创建一个新的顶点 g 。保留两个 f 的最左儿子作为 f 的儿子, 但是使两个 f 的最右儿子改变为 g 的儿子。然后, 通过使 g 成为顶点 f 的父亲, 使得 g 成为顶点 f 的一个兄弟。如果 f 的父亲有两个儿子, 则到此为止。如果 f 的父亲有 3 个儿子, 则必须递归地重复该过程, 直到树中所有的顶点至多有 3 个儿子。如果根有 4 个儿子, 则创建拥有两个新儿子的新根, 每个儿子拥有先前根的 4 个儿子中的两个。

例 4.10 如果插入元素 4 到图 4-27a 的 2-3 树中, 会发现标记为 4 的新叶子应该成为顶点 c 的最左儿子。由于顶点 c 已经有了 3 个儿子, 则创建一个新顶点 c' 。使叶子 4 和 5 成为 c 的儿子, 叶子 6 和 7 成为 c' 的儿子。现在使 c' 成为顶点 a 的儿子。然而, 由于顶点 a 已经有了 3 个儿子, 因此创建另一个顶点 a' 。使顶点 b 和 c 成为旧顶点 a 的儿子, 顶点 c' 和 d 成为新顶点 a' 的儿子。最终, 创建新根 r , 并使 a 和 a' 成为 r 的儿子。得到的树如图 4-28 所示。 □

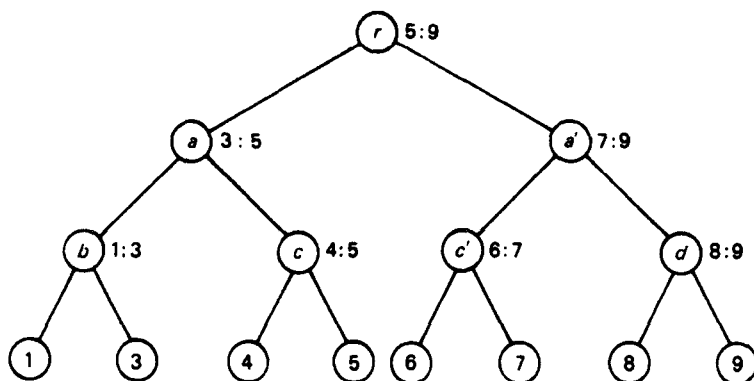


图 4-28 在图 4-27a 的树中插入 4

算法 4.4 插入新元素到一棵 2-3 树中。

输入: 一棵根为 r 的非空 2-3 树 T , 一个不在 T 中的新元素 a 。

输出: 一棵修正后的 2-3 树, 其包含标记为 a 的新叶子。

方法: 假定 T 至少有一个元素。为简化算法描述, 忽略对各顶点更新 L 和 M 值的细节。

1. 如果 T 包含标记为 b 的唯一叶子 l , 则创建一个新根 r' 。创建一个标记为 a 的新叶子 v 。使 l 和 v 成为 r' 的儿子, 如果 $b < a$, 则使 l 成为左儿子; 否则, 使 l 成为右儿子。

2. 如果 T 的顶点超过一个, 设 $f = \text{SEARCH}(a, r)$, 其中 SEARCH 为图 4-29 中的程序。创建一个标记为 a 的新叶子。如果 f 有两个儿子, 标记为 b_1 和 b_2 , 则使 l 成为 f 适当位置的儿子。如果 $a < b_1$, 使 l 成为左儿子; 如果 $b_1 < a < b_2$, 使 l 成为中儿子; 如果 $b_2 < a$, 使 l 成为右儿子。如果 f 包含 3 个儿子, 则使 l 成为 f 适当位置的儿子, 然后调用 $\text{ADDSON}(f)$ 将 f 及其 4 个儿子加入到 T 中。 ADDSON 为图 4-30 中的程序。沿着从 a 到根的路径调整 L 和 M 值, 来解释 a 的出现。[⊙] 其他对 L 和 M 值的显著调整由 ADDSON 来完成(此处略, 留作练习)。□

```

procedure SEARCH( $a, r$ ):
if any son of  $r$  is a leaf then return  $r$ 
else
  begin
    let  $s_i$  be the  $i$ th son of  $r$ ;
    if  $a \leq L[r]$  then return SEARCH( $a, s_1$ )
    else
      if  $r$  has two sons or  $a \leq M[r]$  then return SEARCH( $a, s_2$ )
      else return SEARCH( $a, s_3$ )
  end

```

图 4-29 SEARCH 程序

定理 4.6 算法 4.4 插入一个新元素到含有 n 个叶子的一棵 2-3 树的操作最多消耗 $O(\log n)$ 时间。而且, 算法保持了原有叶子的顺序并保留了 2-3 树的结构。

证明: 使新叶子成为正确顶点的儿子的 SEARCH 过程的调用次数进行直接归纳。注意原有叶子的顺序并没有被打乱。就时间而言, 由引理 4.6 包含 n 片叶子的一棵 2-3 树的高度至多为 $\log n$ 。由于 $\text{ADDSON}(v)$ 仅对 v 的父亲递归地调用自身, 至多可能有 $\log n$ 次递归调用。由于 ADDSON 的每次调用仅要求常数时间, 因而总时间至多为 $O(\log n)$ 。□

```

procedure ADDSON( $v$ ):
begin
  create a new vertex  $v'$ ;
  make the two rightmost sons of  $v$  the left and right sons of  $v'$ ;
  if  $v$  has no father then
    begin
      create a new root  $r$ ;
      make  $v$  the left son and  $v'$  the right son of  $r$ 
    end
  else
    begin
      let  $f$  be the father of  $v$ ;
      make  $v'$  a son of  $f$  immediately to the right of  $v$ ;
      if  $f$  now has four sons then ADDSON( $f$ )
    end
  end

```

图 4-30 ADDSON 程序

能够从 2-3 树中删除元素 a , 所用方法必然与插入一个元素相反。假定元素 a 是叶子 l 的标记。考虑 3 种情形。

情形 1 如果 l 为根, 删除 l (在这种情况下, a 是树中的唯一元素);

情形 2 如果 l 为拥有 3 个儿子的顶点的儿子, 删除 l ;

⊙ 仅需沿着从 a 开始的路径, 直至到达一个顶点, a 不是以该顶点为根的子树中的最大元素。

情形 3 如果 l 为拥有两个儿子 s 和 l 的顶点 f 的儿子, 则有两种可能性:

a) f 为根。删除 l 和 f , 保留剩下的儿子 s 作为根。

b) f 不是根。假定 f 有一个兄弟 g^{\ominus} 位于其左侧。当兄弟位于其右侧时, 可相似处理。如果 g 仅有两个儿子, 则使 s 为 g 的最右儿子, 删除 l , 递归调用删除程序删除 f 。如果 g 有 3 个儿子, 使 g 的最右儿子为 f 的左儿子并从树中删除 l 。

例 4.11 从图 4-28 的 2-3 树中删除元素 4。标记为 4 的叶子是拥有两个儿子的顶点 c 的儿子。因此, 使标记为 5 的叶子成为顶点 b 的最右儿子, 删除标记为 4 的叶子, 然后递归地删除顶点 c 。

顶点 c 为顶点 a 的儿子, c 自己有两个儿子。顶点 a' 为 a 的右兄弟。因此, 由对称性, 使 b 成为 a' 的最左儿子, 删除顶点 c , 然后递归地删除顶点 a 。

顶点 a 是树根的儿子。应用情形 3(a), 保留 a' 作为剩余树的根, 如图 4-31 所示。 □

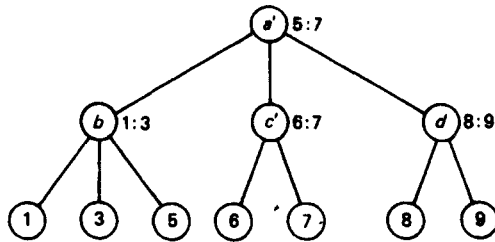


图 4-31 在图 4-28 的树中删除 4 之后

对于 n 个叶子的 2-3 树, 可以在至多 $O(\log n)$ 步内完成该过程的执行。其证明连同删除过程的形式说明一起留作练习。

至此, 已经证明了对一棵含有 n 个叶子的 2-3 树, 至多需要 $O(\log n)$ 步执行 MEMBER、INSERT 或 DELETE 指令。因此, 一棵 2-3 树可用来作为代价为 $O(n \log n)$ 的字典, 因为它能够在至多 $O(n \log n)$ 步内, 处理包含 n 条 MEMBER、INSERT 和 DELETE 的指令序列。

现在考虑 MIN 指令。在一棵 2-3 树中, 最小元素位于最左叶子上, 肯定能够在 $O(\log n)$ 步内被找到。因此, 任何含有 n 条 INSERT、DELETE 和 MIN 指令的指令序列都能够通过使用一棵 2-3 树, 在 $O(n \log n)$ 时间内处理。因此, 已经验证了前面声称的事实: 一棵 2-3 树可以用来实现一个 $O(n \log n)$ 优先队列。其他两个可用来实现 $O(n \log n)$ 优先队列的数据结构是在堆排序 (Heapsort) 中所用到的堆, 以及在习题 4.30 ~ 4.33 中所讨论的 AVL 树。

4.11 可合并堆

本节给出一个数据结构, 可在 $O(n \log n)$ 时间内处理由 n 条 INSERT、DELETE、UNION 和 MIN 指令组成的指令序列。可认为该结构是 3.4 节所讨论的堆的一般化, 用一棵 2-3 树 T 表示元素集 S 。 S 中的每个元素作为 T 的叶子的标记出现, 但是, 不同于前面两节, 叶子不是线性有序的。对于 T 中的每个内部顶点, 赋予标记 $\text{SMALLEST}[v]$, 表示存储在以 v 为根的子树中最小元素的值。在 2-3 树的这一应用中, 不需要 $L[v]$ 和 $M[v]$ 。

从 T 的根开始, 沿着树按如下方式下移, 可找到集合 S 中的最小元素。如果当前位于内部顶点 v , 接下来访问 SMALLEST 值最小的 v 的儿子。如果 T 有 n 片叶子, 则指令 MIN 需要 $O(\log n)$ 步。

在多数应用中, 无论何时需要从 S 中删除一个元素, 该元素总是最小元素。然而, 如果要

\ominus 有共同父亲的两个顶点称为兄弟。

从 S 中删除任意一个元素, 则必须能够找到包含该元素的叶子。在元素可用整数 $1, 2, \dots, n$ 表示的应用中, 可以直接索引到叶子。如果元素是任意的, 可以使用一个辅助 2-3 字典, 其叶子包含指向 T 中叶子的指针。通过该字典, 可以在 $O(\log n)$ 步内访问 T 中的任意一片叶子。每当执行一条 INSERT 指令的时候, 必须更新字典, 这至多需要 $O(\log n)$ 步。

一旦从 T 中删除了一片叶子, 则必须为 l 的每个特定的祖先 v 重新计算 SMALLEST 值。SMALLEST[v] 的新值将是 v 的两或三个儿子 s 中的 SMALLEST[s] 的最小值。如果总是自底向上重新计算, 通过对重新计算的次数进行归纳, 可以证明每次计算都产生 SMALLEST 的正确答案。由于在被删除叶子的祖先上仅有 SMALLEST 值发生改变, 所以可在 $O(\log n)$ 步内处理完 DELETE 指令。

现在考虑 UNION 指令。每个集合都用一棵不同的 2-3 树表示。为了合并两个集合 S_1 和 S_2 , 调用图 4-32 的程序 IMPLANT(T_1, T_2), 其中 T_1 和 T_2 为表示 S_1 和 S_2 的 2-3 树^①。

假定 T_1 的高度 h_1 大于等于 T_2 的高度 h_2 。在 T_1 的最右路径上, IMPLANT 找到高度为 h_2 的顶点 v , 使 T_2 的根为 v 的最右兄弟。如果 v 的父亲 f 现在有了 4 个儿子, 则 IMPLANT 调用程序 ADDSON(f)。在 IMPLANT 操作期间, 后代发生改变的顶点的 SMALLEST 值可以用 DELETE 操作中一样的方式更新。

IMPLANT 程序可在 $O(h_1 - h_2)$ 时间内合并 T_1 和 T_2 到一棵 2-3 树, 证明留作练习。当我们计算更新 L 和 M 值的时间的时候, 程序 IMPLANT 可能需要用 $O(\text{MAX}(\log \|S_1\|, \log \|S_2\|))$ 时间。

```

procedure IMPLANT( $T_1, T_2$ ):
  if HEIGHT( $T_1$ ) = HEIGHT( $T_2$ ) then
    begin
      create a new root  $r$ ;
      make ROOT[ $T_1$ ] and ROOT[ $T_2$ ] the left and right sons of  $r$ 
    end
  else
    wlg assume HEIGHT( $T_1$ ) > HEIGHT( $T_2$ ) otherwise
    interchange  $T_1$  and  $T_2$  and interchange "left" and "right" in
    begin
      let  $v$  be the vertex on the rightmost path of  $T_1$  such that
        DEPTH( $v$ ) = HEIGHT( $T_1$ ) - HEIGHT( $T_2$ );
      let  $f$  be the father of  $v$ ;
      make ROOT[ $T_2$ ] a son of  $f$  immediately to the right of  $v$ ;
      if  $f$  now has four sons then ADDSON( $f$ )①
    end
  end

```

- ① 如果使 ADDSON(f) 产生的新顶点取值 $< M$, 必须先沿着到最右边叶结点的路径求出 v 的最大后代

图 4-32 IMPLANT 程序

下面考虑一个应用, 其中 UNION、MIN 和 DELETE 的操作是随意出现的。

例 4.12 考虑一个算法, 其查找例 4-1 中图的最小代价生成树。顶点存放到逐渐变大的集合中, 通过由最小代价生成树选出的边来连接每个集合中的成员。为生成树找出新边就是考察边(首先是最短的), 观察它们是否与尚未连接的顶点相连接。

另一个策略是对每个顶点集 V_i , 维护同 V_i 中的某顶点相关联的所有未考虑的边的集合 E_i 。如果选择一条先前未考虑过的边 e , 该边与相对较小集合 V_i 中的某个顶点关联, 则有很大可能 e 的另一端不在 V_i 中, 可以将 e 加入到生成树中。如果 e 是同 V_i 中的顶点相关联的且未考虑的最小代价边, 那么可以证明, 添加 e 到生成树将导出一棵最小代价生成树。

① 在此, “左”“右”之间的差异并不重要, 但是这种差异在下一节讨论可连接队列时需要考虑。

使用刚才所描述的树连接算法, 合并路径左侧的树和仅包含 a 的树。同样地, 合并路径右侧的树。图 4-34 给出的程序 DIVIDE 包含了具体细节。

```

procedure DIVIDE( $a, T$ ):
begin
  on the path from  $\text{ROOT}[T]$  to the leaf labeled  $a$  remove all vertices except the leaf;
  comment At this point  $T$  has been divided into two forests—the left forest, which consists of all trees with leaves to the left of and including the leaf labeled  $a$ , and the right forest, which consists of all trees with leaves to the right of  $a$ ;
  while there is more than one tree in the left forest do
    begin
      let  $T'$  and  $T''$  be the two rightmost trees in the left forest;
      IMPLANT( $T', T''$ )①
    end;
  while there is more than one tree in the right forest do
    begin
      let  $T'$  and  $T''$  be the two leftmost trees in the right forest;
      IMPLANT( $T', T''$ )
    end
  end

```

① IMPLANT(T', T'')的结果可考虑为左森林中所留下的。类似地, 当应用于右森林中的树时, IMPLANT 的结果是右森林中的一棵树。

图 4-34 分离一棵 2-3 树的程序

定理 4.7 程序 DIVIDE 以叶子 a 划分 2-3 树 T , 使 a 左侧的所有叶子与 a 本身位于一棵 2-3 树, 且 a 右侧的所有叶子位于第 2 棵 2-3 树。该程序耗时 $O(\text{HEIGHT}(T))$, 而叶子的顺序不变。

证明: 依照程序 IMPLANT 的性质, 树适当地重新装配如下。通过观察下述事实可得到时间开销的结果。首先, 对应任意给定的高度, 除了高度为 0 可以对应 3 棵树外, 其他则至多存在两棵树。当合并两棵树时, 结果树的高度至多比两棵原有树高度的最大值大 1。在结果树的高度比每一棵原始树的高度大 1 的情况下, 其根的度为 2。因此, 如果合并高度为 h 的 3 棵树, 则结果树的高度至多为 $h+1$ 。因此, 在重组过程的每个阶段, 相同高度的树至多存在 3 棵。

由于合并两棵不同高度的树所需要的时间同它们高度的差异成比例, 且合并两棵相同高度的树所需要的时间为常量, 所以重新合并所有树的时间, 与树的数目加上任意两棵树高度的最大差异值的和成比例。因此, 所花费的总时间同原始树的高度同阶。□

通过观察可以看到, 利用可连接队列, 能够在 $O(\text{MAX}(\log |S_1|, \log |S_2|))$ 时间内将序列 S_2 插入到序列 S_1 中的一对元素之间。如果 $S_2 = b_1, b_2, \dots, b_n$, $S_1 = a_1, a_2, \dots, a_m$, S_2 要插入到元素 a_i 和 a_{i+1} 之间, 那么可以使用指令 $\text{SPLIT}(a_i, S_1)$, 基于 a_i 将 S_1 划分为两个序列 $S_1' = a_1, \dots, a_i$ 和 $S_1'' = a_{i+1}, \dots, a_m$ 。然后, 用 $\text{CONCATENATE}(S_1', S_2)$ 得到序列 $S_3 = a_1, \dots, a_i, b_1, \dots, b_n$, 最终通过 $\text{CONCATENATE}(S_3, S_1'')$ 得到需要的序列。

4.13 划分

现在考虑一种特殊的集合分离, 称为划分(partitioning)。划分一个集合的问题经常出现, 在此给出的解决方案是启发式的。假定给定一个集合 S 以及 S 的初始划分 π , 将 S 划分成不相交的块 $\{B_1, B_2, \dots, B_p\}$ 。也给定关于 S 的函数 f 。任务是找到 S 的最粗粒度(coarsest)(包含的块最少)的划分, 称为 $\pi' = \{E_1, E_2, \dots, E_q\}$, 使得:

1. π' 与 π 一致(即每个 E_i 都是 B_j 的子集);
2. E_i 中的 a 和 b 意味着 $f(a)$ 和 $f(b)$ 处于某个 E_j 中。

将 π' 称为与 π 和 f 兼容的 S 的最粗粒度划分。

显而易见的解决方案是使用下述方法不断改进原始划分中的块。设 B_i 为一块。对 B_i 中的每个 a ，验证 $f(a)$ 。然后划分 B_i ，当且仅当 $f(a)$ 和 $f(b)$ 都处于某块 B_j 中时，元素 a 和 b 放入到相同的块中。重复该过程，直到不可能进一步改进为止。由于每次改进都需要 $O(n)$ 时间，且可能存在 $O(n)$ 次改进，这种方法将得到一个 $O(n^2)$ 的算法。事实上，这种方法确实可能需要平方数量级的计算步数，例 4-13 可说明该结论。

例 4.13 设 $S = \{1, 2, \dots, n\}$ ， $B_1 = \{1, 2, \dots, n-1\}$ 和 $B_2 = \{n\}$ 为初始划分。设 f 为 S 的函数，对于 $1 \leq i < n$ ， $f(i) = i+1$ 且 $f(n) = n$ 。第一次重复过程时，将 B_1 划分为 $\{1, 2, \dots, n-2\}$ 和 $\{n-1\}$ 。由于必须验证 B_1 中的每个元素，这次重复需要 $n-1$ 步。下一次重复时，将 $\{1, 2, \dots, n-2\}$ 划分为 $\{1, 2, \dots, n-3\}$ 和 $\{n-2\}$ 。继续按这种方法处理，总共需要 $n-2$ 次重复，第 i 次重复花费 $n-i$ 步。因此，总共需要的步数为

$$\sum_{i=1}^{n-2} (n-i) = \frac{n(n-1)}{2} - 1$$

对于 $1 \leq i \leq n$ ，最终的划分得到 $E_i = \{i\}$ 。

该方法的困难在于，即使仅从块中删除一个元素，改进一个块也可能要 $O(n)$ 步。这里将设计一个划分算法，其中改进一个块为两个子块所需要的时间同较小的子块成比例。该方法将得到一个 $O(n \log n)$ 的算法。

对于每个 $B \subseteq S$ ，令 $f^{-1}(B) = \{b \mid f(b) \in B\}$ 。对于 $a \in B_i$ ，不是通过 $f(a)$ 的值划分块 B_i ，而是相对于 B_i ，划分包含至少一个 $f^{-1}(B_i)$ 中元素和一个不在 $f^{-1}(B_i)$ 中元素的块 B_j 。也就是说，每个这样的 B_j 都划分为集合 $\{b \mid b \in B_j \text{ 且 } f(b) \in B_i\}$ 和 $\{b \mid b \in B_j \text{ 且 } f(b) \notin B_i\}$ 。

一旦关于 B_i 进行了划分，就不需要再关于 B_i 进行划分，除非 B_i 本身被分裂。如果初始对于每个元素 $b \in B_j$ ， $f(b) \in B_i$ ， B_i 分裂为 B_i' 和 B_i'' ，则可以关于 B_i' 或 B_i'' 对 B_j 进行划分。由于 $\{b \mid b \in B_j \text{ 且 } f(b) \in B_i'\}$ 等于 $B_j - \{b \mid b \in B_j \text{ 且 } f(b) \in B_i''\}$ ，我们将得到相同的结果。

对于关于 B_i' 或 B_i'' 的划分已经有了选择，为此划分更为容易的那个。用 $f^{-1}(B_i')$ 和 $f^{-1}(B_i'')$ 中更小的那个进行划分。图 4-35 给出了算法。

```

begin
1.  WAITING ← {1, 2, ..., p};
2.  q ← p;
3.  while WAITING not empty do
      begin
4.      select and delete any integer i from WAITING;
5.      INVERSE ← f-1(B[i]);
6.      for each j such that B[j] ∩ INVERSE ≠ ∅ and
          B[j] ⊄ INVERSE do
              begin
7.                  q ← q + 1;
8.                  create a new block B[q];
9.                  B[q] ← B[j] ∩ INVERSE;
10.                 B[j] ← B[j] - B[q];
11.                 if j is in WAITING then add q to WAITING
              else
12.                  if ||B[j]|| ≤ ||B[q]|| then
13.                      add j to WAITING
14.                  else add q to WAITING
              end
      end
end

```

图 4-35 划分算法

算法 4.5 划分。

输入：含有 n 个元素的集合 S ，初始划分 $\pi = \{B[1], \dots, B[p]\}$ 以及函数 $f: S \rightarrow S$ 。

输出：划分 $\pi' = \{B[1], B[2], \dots, B[q]\}$ ，使得 π' 为与 π 和 f 兼容的 S 的最粗粒度划分。

方法：将图 4-35 中的程序应用于 π 。该程序忽略了某些重要的实现细节。稍后，在分析运行时间的时候，将讨论那些细节。 \square

现在开始分析算法 4.5，假定其正确地划分了 S 。设 π 为集合 S 的任意划分， f 为 S 的一个函数。如果对 π 中的每块 B ，要么 $B \subseteq f^{-1}(T)$ ，要么 $B \cap f^{-1}(T) = \emptyset$ ，则称集合 $T \subseteq S$ 对 π 是安全的。例如，在图 4-35 中第 9 和 10 行确保对于最终的划分 $B[i]$ 是安全的，这是由于，如果对于某块 B ， $B \cap f^{-1}(B[i]) \neq \emptyset$ ，则或者 $B \subseteq \text{INVERSE}$ ，直接可得 $B \subseteq f^{-1}(B[i])$ ，或者在第 9 和 10 行， B 分成两个块，一块为 $f^{-1}(B[i])$ 的子集，另一块同该集合不相交。

算法 4.5 正确性证明的部分工作在于，证明最后的划分粒度不会太粗。即必须证明下列引理。

引理 4.7 算法 4.5 终止后，在最终划分 π' 中的每块 B 对 π' 是安全的。

证明：实际上，对于每块 $B[l]$ ，所要证明的是：

在每次执行图 4-35 中第 4~14 行的循环之后，如果 l 不在 WAITING 中，且 $l \leq q$ ，那么存在列表 q_1, q_2, \dots, q_k (可能为空)，使得对 $1 \leq i \leq k$ ，每个 q_i 都在 WAITING 中，且 $B[l] \cup B[q_1] \cup \dots \cup B[q_k]$ 对于当前的划分是安全的 (4-6)

直观地，当从 WAITING 中删除块 $B[l]$ 时，第 6~14 行使得 $B[l]$ 对于第 14 行之后的划分是安全的。 $B[l]$ 保持安全性，直到被划分。当 $B[l]$ 被划分时，一个称为 $B[q]$ 的子块被放到 WAITING 中。另一个子块仍然称为 $B[l]$ 。显然，这两个子块的并 $B[l] \cup B[q]$ 是安全的，因为它是原先的 $B[l]$ 。进一步的划分产生了块 $B[q_1], \dots, B[q_k]$ ，使得 q_1, q_2, \dots, q_k 处于 WAITING 中，且 $R = B[l] \cup B[q_1] \cup \dots \cup B[q_k]$ 是安全的。当从 WAITING 中删除某个 q_i ， $1 \leq i \leq k$ 时，第 6~14 行再次使 $B[q_i]$ 和 $R - B[q_i]$ 是安全的。下面将更精确地给出这些思想。

通过对执行第 4~14 行的次数进行归纳可证明式 (4-6) 对所有的 l 都成立。当算法终止的时候，WAITING 为空，因此式 (4-6) 意味着在最终划分 π' 中的每块对 π' 是安全的。

归纳基础的执行次数为 0，由于对所有的 $1 \leq l \leq q = p$ ， l 在 WAITING 中，于是式 (4-6) 显然。

对于归纳步，假定在完成第 14 行后， l 不在 WAITING 中。如果之前的执行之后 l 在 WAITING 中，那么 l 具有在第 4 行中定义的值 i 。易证，第 6~14 行中的循环使 $B[i]$ 对于第 14 行之后的实际划分是安全的。在“安全”性的定义之后，我们曾证明过该事实。

如果在执行第 14 行之前的各语句执行之后， l 不在 WAITING 中，则由归纳假设，存在列表 q_1, q_2, \dots, q_k ，使得对前一阶段中的 l ，满足式 (4-6)。也可以保证，在第 4 行， $i \neq l$ 。

情形 1 i 不在 $L = \{q_1, q_2, \dots, q_k\}$ 之中。在第 9~10 行可以分离为几个块。对于每 $B[q_r]$ ， $1 \leq r \leq k$ (即 $j = q_r$)，将在第 8 行所创建的块的下标加入 L 。由第 11 行， L 仍然将仅由 WAITING 中的下标组成。如果 $B[l]$ 本身没有分解，则 $B[l]$ 以及 L 中块的集合仍然组成一个对当前划分安全的集合，满足式 (4-6)。如果 $B[l]$ 被分解，也必须当 $j = l$ 时把第 8 行选择的索引 q 加入 L 。因此， $B[l] \cup \bigcup_{r \in L} B[r]$ 对于当前的划分将是安全的。

情形 2 i 在 $L = \{q_1, q_2, \dots, q_k\}$ 之中。不失一般性，假定 $i = q_1$ 。证明过程大体上类似于情形 1，不过在执行第 4~14 行的当前循环之后， q_1 可能不在 WAITING 中。然而，我们知道，当前划分的每块 B 要么是 $f^{-1}(B[q_1])$ 的一个子集，要么同它不相交。令 $T = B[l] \cup \bigcup_{r \in L} B[r]$ ，其中 L 已按情形 1 修正。如果 $B \subseteq f^{-1}(B[q_1])$ ，那么肯定有 $B \cap f^{-1}(T) = \emptyset$ 。如果 $B \cap f^{-1}(B[q_1]) = \emptyset$ ，那么，或者有 $B \cap f^{-1}(T) = \emptyset$ 或者有 $B \subseteq f^{-1}(T)$ ，其证明同情形 1 相似。

最终, 当算法 4.5 终止时, WAITING 必然为空。因此, 式(4-6)意味着, 对于每个 l , $B[l]$ 对于最终的划分是安全的。□

定理 4.8 算法 4.5 正确地计算了与 π 和 f 可兼容的 S 的最粗划分。

证明: 引理 4.7 表明, 算法 4.5 的输出 π' 同 π 和 f 是可兼容的。需要证明, π' 的粒度是尽可能粗的。可简单地对由第 9~10 行分离得到的块数进行归纳, 以证明算法每次进行这样的分离对于可兼容性是必要的。归纳证明留作练习。□

为了证明算法 4.5 的运行时间为 $O(n \log n)$, $n = \|S\|$, 现在必须详细考虑算法的实现。证明时间开销的关键在于, 要证明第 6~14 行的循环怎么能够在同 $\|INVERSE\|$ 成比例的时间内执行。第一个问题是, 在第 6 行怎样有效地找到 j 的合适的集合。我们需要一个数组 INBLOCK, 使得 INBLOCK[l] 为包含 l 的块的索引。可以在 $O(n)$ 步内初始化 INBLOCK, 并在第 9 行之后更新, 且其耗时不超过在该行创建列表 $B[q]$ 所消耗的时间。因此, 由于仅关心时间复杂度大小的阶数, 忽略 INBLOCK 的处理是合理的。

使用 INBLOCK, 很容易在 $O(\|INVERSE\|)$ 步内构造一个在第 6 行中所需要的 j 的列表 JLIST。对于 INVERSE 中的每个元素 a , 如果包含 a 的块的索引不在 JLIST 中, 则加入该索引到 JLIST。对每个 j , 计算在 INVERSE 中也在 $B[j]$ 中的元素的数目。如果该数目达到 $\|B[j]\|$, 则 $B[j] \subseteq INVERSE$, 则从 JLIST 中删除 j 。

使用 INBLOCK, 对于 JLIST 中的每个 j , 也可以构建一个整数列表 INTERSECTION[j], 以存放 $B[j]$ 和 INVERSE 的交集。为了从 $B[j]$ 中快速删除在 INTERSECTION[j] 中的元素, 并将其加入 $B[q]$, 必须以双链表的形式维护列表 $B[l]$, $1 \leq l \leq q$, 即, 指针指向后继和前驱。

第 9 和 10 行需要 $O(\|B[q]\|)$ 步。对于 for 循环的给定执行, 找到合适的 j 并执行 7~10 行所耗费的总时间为 $O(\|INVERSE\|)$ 。对第 12~14 行进行验证, 易发现如果操作适当, 相对于总时间 $O(\|INVERSE\|)$, 它需要 $O(\|B[q]\|)$ 时间。

现在仅需要考虑第 11 行。为快速地断定 j 是否在列表 WAITING 中, 创建另一个数组 INWAITING[j]。可在 $O(n)$ 步内初始化 INWAITING, 并且在第 11~14 行中毫不费力地对其进行维护。因此, 有下面的引理。

引理 4.8 在图 4-35 中第 6~14 行的 for 循环可以在 $O(\|INVERSE\|)$ 步内实现。

证明: 由上面所述。□

定理 4.9 可在 $O(n \log n)$ 时间内实现算法 4.5。

证明: 考虑一个整数 s 所在的块不在 WAITING 中, 能够在放入到 WAITING 中的块中找到其对应的块。在第 1 行, 这可能出现一次。在第 11 行, 即使 $s \in B[q]$, 这不可能出现, 因为 s 先前已经处在 $B[j]$ 中, 且 j 已经在 WAITING 中。如果这在第 13 或 14 行出现, 则 s 处于一块, 相对于包含 s 的集合索引放入 WAITING 中时已包含 s 的块, 该块的大小不到后者的一半。可以得到结论, 包含 s 的集合索引放入到 WAITING 中的次数不会多于 $1 + \log n$ 次。结果, s 处于第 4 行所选块 i 中的次数不可能超过 $1 + \log n$ 次。

假定第 6~14 行中的循环每次执行的代价成比例于 $\|f^{-1}(s)\|$, 由 $B[i]$ 中的元素 s 承担。存在一个常数 c , 对于循环的执行, s 所承担的代价不超过 $c\|f^{-1}(s)\|$ 。但是, 前面已经证明, s 在选择 $B[i]$ 中出现的次数不会多于 $O(\log n)$ 次, 因此它所负担的总代价为 $O(\|f^{-1}(s)\| \times \log n)$ 。由于 $\sum_{s \in S} \|f^{-1}(s)\|$ 必定为 n , for 循环所有执行的总代价为 $O(n \log n)$ 。易见, 算法 4.5 剩余部分的代价为 $O(n)$, 因此上述定理得证。□

算法 4.5 有几种应用。一个重要的应用是在有穷状态机中, 对状态数进行最小化。给定一个有穷自动机 $M = (S, I, \delta, s_0, F)$, 要找到有最少状态数的与 M 等价的自动机 M' 。对于每个状态 s 和输入符号 a , $\delta(s, a)$ 表示 M 的下一状态。起初, M 的状态可以划分为接收状态集 F 和非

接收状态集 $S-F$ 。在 M 中最小化状态数的问题等价于找出 S 的最粗划分 π' ，该划分与初始划分 $\{F, S-F\}$ 一致，如果状态 s 和 t 位于 π' 的一块中，那么，对于每个输入符号 a ，状态 $\delta(s, a)$ 和 $\delta(t, a)$ 也在 M' 的一块中。

该问题同算法 4.5 中问题的唯一差异在于， δ 是从 $S \times I$ 到 S 的映射，而不只是从 S 到 S 的映射。然而，可以将 δ 看作是 S 的函数的集合 $\{\delta_{a_1}, \delta_{a_2}, \dots, \delta_{a_l}\}$ ，其中，每个 δ_{a_i} 都是 δ 对输入符号 a_i 的约束。

通过将元素对 (i, δ_{a_i}) 放入集合 WAITING，可以很容易地将算法 4.5 修改为处理这种更为一般的问题。每个元素对 (i, δ_{a_i}) 包含划分中一块的索引 i ，加上该划分的函数 δ_{a_i} 。起初， $\text{WAITING} = \{(i, \delta_{a_i}) \mid i=1 \text{ 或 } 2, \text{ 且 } a_i \in I\}$ ，因为初始划分 $\{F, S-F\}$ 包含两块。任何时候，块 $B[j]$ 分离成 $B[j]$ 和 $B[q]$ ，每个可能的函数 δ_{a_i} 都以 j 和 q 来配对。余下的细节留作练习。

4.14 本章小结

图 4-36 总结了本章所讨论的各种数据结构，它们能够处理的指令类型以及所作的元素所取自的全集的大小和元素性质的假设。

数据结构	全集类型	允许的指令	对大小为 n 的集合处理 n 条指令的时间	
			期望时间	最坏情况下时间
1. 散列表	散列函数可计算的任意集合	MEMBER, INSERT, DELETE	$O(n)$	$O(n^2)$
2. 二叉查找树	任意有序集	MEMBER, INSERT, DELETE, MIN	$O(n \log n)$	$O(n^2)$
3. 算法 4.3 的树结构	整数 $1 \sim n$	MEMBER, INSERT, DELETE, UNION, FIND	至多 $O(nG(n))$	至多 $O(nG(n))$
4. 叶子无序的 2-3 树	任意有序集	MEMBER, INSERT, DELETE, UNION, FIND, MIN	$O(n \log n)$	$O(n \log n)$
5. 叶子有序的 2-3 树	任意有序集	MEMBER, INSERT, DELETE, FIND, SPLIT, MIN, CONCATENATE	$O(n \log n)$	$O(n \log n)$

图 4-36 数据结构的性质总结

习题

- 4.1 给出 4.1 节中 7 个基本操作的某个子集，使该子集足以对 n 个元素组成的任一序列进行排序。解释执行从你的子集中选出的 n 条指令流的复杂度。
- 4.2 假定元素是字母串，并且对于大小 $m=5$ 的表，使用如下散列函数：给字母赋值， A 的值为 1， B 的值为 2，依此类推。将这些字母“值”相加，将和的值除以 5 并取余数。在插入串 DASHER, DANCER, PRANCER, VIXEN, COMET, CUPID, DONNER, BLITZEN 后，给出散列表和列表的内容。
- 4.3 将习题 4.2 中的 8 个串插入一棵二叉查找树。如果要验证 RUDOLPH 是否是集合中的成员，被访问的顶点序列是什么？
- 4.4 证明：如果能够等可能地看到全集中的所有元素，那么图 4-3 中的 SEARCH 程序给出的可能期望查找时间是最小的。
- 4.5 找出一棵最优二叉查找树，使元素 a, b, \dots, h 出现的概率依次为 0.1, 0.2, 0.05, 0.1, 0.3, 0.05, 0.15, 0.05，且所有其他元素的概率都为 0。
- **4.6 证明在算法 4.2 中，若将图 4-9 第 8 行中寻找 m 的范围限定在从位置 $r_{i,j-1}$ 到位置 $r_{i+1,j}$ ，仍

然可以找到最大值。

- *4.7 用习题 4.6 的结论修改算法 4.2, 使其运行时间为 $O(n^2)$ 。
- 4.8 通过证明图 4-10 树构造算法的正确性, 以完成定理 4.2 的证明。
- 4.9 完成定理 4.3 的证明。
- *4.10 在从 1 到 r 的 n 个外部名组成的集合与包含整数 $1 \sim n$ 的内部名组成的集合之间, 构造一个接口使其可以相互转换。该接口必须满足: 能够在两个方向进行转换。假定 $r > n$ 。
 - a) 设计接口, 使其具有良好的期望时间性能;
 - b) 设计接口, 使其具有良好的最坏情况性能。
- *4.11 给出表示由从 1 到 n 的整数组成的子集 S 的有效的数据结构。欲对该集合执行的操作有:
 - 1) 从集合中任意选择一个整数, 并将其删除;
 - 2) 将整数 i 加入集合。
 必须提供一种机制, 在 S 已经包含 i 的情况下, 忽略加入整数 i 到 S 的请求。该数据结构必须满足: 选择并删除一个元素的时间和插入一个元素的时间是不依赖于 $\|S\|$ 的常量。
- 4.12 用算法 4.3 执行根据下列程序生成的 UNION 和 FIND 指令序列, 给出由此得到的树。假定, 对于 $1 \leq i \leq 16$, 初始时集合 i 为 $\{i\}$ 。

```

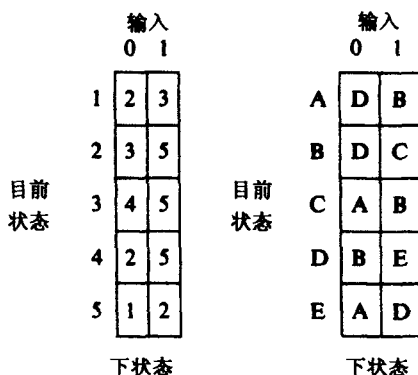
begin
  for  $i \leftarrow 1$  step 2 until 15 do UNION( $i, i+1, i$ );
  for  $i \leftarrow 1$  step 4 until 13 do UNION( $i, i+2, i$ );
  UNION(1, 5, 1);
  UNION(9, 13, 9);
  UNION(1, 9, 1);
  for  $i \leftarrow 1$  step 4 until 13 do FIND( $i$ )
end

```

- 4.13 设 σ 为 UNION 和 FIND 指令组成的序列, 其中, 所有的 UNION 都出现在 FIND 指令之前。证明: 算法 4.3 执行 σ 时, 执行时间同 σ 的长度成比例。
- 4.14 S_0 树为一棵包含单个顶点的树。对于 $i > 0$, 通过使一棵 S_{i-1} 树的根成为另一棵 S_{i-1} 树的根的儿子, 得到一棵 S_i 树。证明如下结论:
 - a) 一棵 S_n 树包含 $\binom{n}{h}$ 个高为 h 的顶点;
 - b) 对于 $m \leq n$, 通过将 S_m 树中的每个顶点替换为一棵 S_{n-m} 树, 可以由一棵 S_m 树得到一棵 S_n 树。此时, 各个顶点的儿子变成替换进来的 S_{n-m} 树的根的儿子;
 - c) 一棵 S_n 树包含一个有 n 个儿子的顶点, 这些儿子分别为一棵 S_0 树, 一棵 S_1 树, \dots , 一棵 S_{n-1} 树的根。
- 4.15 为不相交集合并问题给出一个算法, 使得包含较少顶点(可任意打乱此约束)的树的根成为顶点数较大的树的根的儿子, 但不使用路径压缩。证明: $O(n \log n)$ 的上界无法再得到改进。即, 对某个常数 c 和任意大的 n , 存在 UNION 和 FIND 指令组成的序列, 其执行需要 $cn \log n$ 步。
- **4.16 设 $T(n)$ 为执行 n 条 UNION 和 FIND 指令的最坏时间复杂度。对于 FIND 指令, 使用 4.7 节中用了路径压缩的树结构, 但是, 通过使 A 的根成为 B 的根的儿子, 并不考虑集合的大小, 来执行 UNION(A, B, C)。证明: 对于某个常数 $k_1 > 0$, $T(n) \geq k_1 n \log n$ 。
- **4.17 证明: 对于某个常数 k_2 , $T(n) \leq k_2 n \log n$, 其中 $T(n)$ 的定义同习题 4.16。
- *4.18 证明可以在 $O(nG(n))$ 时间内, 基于整数 $1, 2, \dots, n$, 执行 n 条 UNION、FIND、MEMBER、INSERT、DELETE 指令组成的序列。假定 DELETE(i, S) 使 i 成为一个新集合 $\{i\}$ 的成员, 可以给该集合一个(任意的)名字。最终这一集合合并到另一个集合。也假定不会

有元素同时作为多个集合的成员。

- 4.19 推广离线 MIN 问题, 处理形如 $\text{MIN}(i)$ 的指令, 该指令确定该 MIN 指令左边的未被前一个 MIN 指令所识别的所有小于 i 的整数。
- 4.20 研究离线 MIN 问题的数据结构, 请用数组表示树, 并基于数组编写程序, 不使用对不相交集合进行合并算法的高级命令。
- 4.21 按如下方式推广离线 MIN 问题。设 T 为一棵树, 包含 n 个顶点。每个顶点同 $1 \sim n$ 之间的一个整数相关联。关联某些顶点是 EXTRACT_MIN 指令。使用后序遍历该树。当碰到顶点 v 上的 EXTRACT_MIN 指令时, 确定并删除以 v 为根的子树上先前没有删除的最小整数 (除了顶点 v 上的整数)。为该过程给出一个离线的 $O(nG(n))$ 算法。
- **4.22 对于 UNION_FIND 问题, 设计一个 $O(n \log \log n)$ 的解决方案, 使用树数据结构, 此类树中的叶子到根的距离为 2。将根的度限定在 $1 \sim n / \log n$ 之间, 根的每个儿子的度则限定为 1 到 $\log n$ 之间。如何修改算法, 使其具有 $O(n \log \log \log n)$ 的时间上界? 通过推广这种方法, 可得到的最优时间上界是什么?
- 4.23 在 4.8 节应用 2 所提到的符号表中, 怎样从初始状态开始跟踪置换过程 [提示: 利用深度测定中所使用的技术]。
- 4.24 为 4.8 节应用 2 中提及的包含有权重计算的 UNION 和 FIND 原语编写一个程序。
- 4.25 用图 4-25 中的算法, 测试下面有穷自动机的等价性。初始状态分别为 1 和 A, 终止状态集分别为 $\{5\}$ 和 $\{C, E\}$ 。



- 4.26 编写完整的在 2-3 树中执行下列指令的程序。
- DELETE ;
 - UNION ;
 - MEMBER (假定叶子有序, 且每个顶点以其最高的叶子作为标记);
 - SPLIT (假定给定了将进行分离的位置的叶子, 且叶子有序)
- 4.27 编写一个完整的插入一个新的顶点到一棵 2-3 树中程序, 假定叶子有序。
- 4.28 为原语 MEMBER 、 INSERT 、 DELETE 、 MIN 、 UNION 和 FIND 编写程序, 其中, 用带有 4.11 节 SMALLEST 标记的 2-3 树。假定全集为 $\{1, 2, \dots, n\}$ 。
- 4.29 考虑可合并堆 (INSERT 、 DELETE 、 UNION 和 MIN) 的 2-3 树实现。假定提取元素的全集足够大。描述怎样在每条 FIND 指令用 $O(\log n)$ 步来实现 FIND , 其中, n 为所包含的所有堆中的元素总数。

定义 AVL 树^①是一棵二叉查找树，在每个顶点 v 上， v 的左右子树的高度相差不大于 1。如果少一棵子树，则认为“高度”-1。

例 4.14 图 4-37 中的树不是一棵 AVL 树，因为标记为 * 的顶点有一棵高为 2 的左子树，且有高度为 0 的右子树。然而，在图 4-37 中，其他的每个顶点都满足 AVL 条件。 □

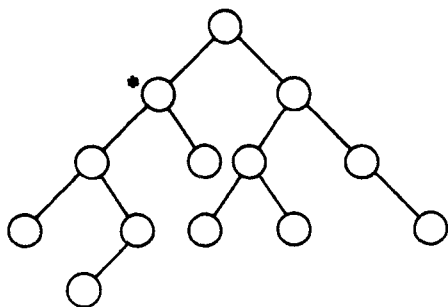


图 4-37 一棵非 AVL 树

- 4.30 证明：一棵高为 h 的 AVL 树至多有 $2^{h+1} - 1$ 个顶点，且至少有

$$\frac{5+2\sqrt{5}}{5}\left(\frac{1+\sqrt{5}}{2}\right)^h + \frac{5-2\sqrt{5}}{5}\left(\frac{1-\sqrt{5}}{2}\right)^h - 1$$

个顶点。

- *4.31 设 T 为一棵包含 n 个顶点的二叉查找树，且其为一棵 AVL 树。为指令 INSERT 和 DELETE 编写 $O(\log n)$ 的算法，使该树仍然为一棵 AVL 树。可以假定，在顶点上能够找到每个顶点的高度，且该高度信息被自动更新。
- *4.32 编写分割一棵 AVL 树和连接两棵 AVL 树的算法。算法的执行时间应同树的高度成比例。
- *4.33 用一棵 AVL 树作为算法的基，对由整数 $1 \sim n$ 组成的集合执行 MIN、UNION 和 DELETE，每个操作使用 $O(\log n)$ 步。

定义 在一棵二叉树中顶点 v 的平衡度 (balance) 为 $(1+L)/(2+L+R)$, 其中, L 和 R 为 v 的左右子树中的顶点数。如果每个顶点的平衡度在 $\alpha \sim 1-\alpha$ 之间, 则称一棵二叉查找树为 α 平衡的。

例 4.15 在图 4-37 中, 标记为 * 的顶点的平衡度为 $5/7$ 。没有其他顶点的平衡度值的偏差大于 $1/2$, 因此, 图 4-37 中的树为 $2/7$ 平衡。 \square

- *4.34 对于一棵高为 h 的 α 平衡树, 给出其顶点数的上下界。
- **4.35 对于平衡度为 α 的树, 重作习题 4.31 ~ 4.33, 其中, $\alpha \leq 1/3$ 。
- *4.36 如果 $\alpha > 1/3$, 能重作习题 4.31 ~ 4.33 吗?
- 4.37 通过叶子上的数据设计一棵平衡树, 通过将子树高度的差异保持在固定常数内, 来得到平衡。
- *4.38 给定一棵 n 个顶点的树以及包含 n 对顶点的列表, 编写一个 $O(nG(n))$ 的算法, 为每个顶点对 (v, w) , 确定为 v 和 w 共同祖先的顶点, 且该顶点为 v 和 w 所有这类共同祖先中最近的。
- *4.39 二叉树的外部路径长度为其深度的所有叶子之和。二叉树的内部路径长度为所有顶点的深度之和。如果每个顶点都有两个儿子, 或者没有儿子, 那么, 外部与内部路径长度之间有什么关系?
- *4.40 设计一个数据结构来实现队列, 以使得下列操作能够在线执行。
 - a) ENQUEUE(i, A): 将整数 i 添加到队列 A 中。允许相同的整数有不同的实例。
 - b) DEQUEUE(A): 从队列 A 中提取出其中最长的元素。
 - c) MERGE(A, B, C): 合并队列 A 和 B , 将得到的队列称为 C 。如果元素在队列 A 或 B 中出现, 则它们认为会出现在合并后的队列中。注意到, 在 A 和 B 中, 同样的整数可能会出现多次, 每一次出现都认为是一个可区分的元素。

对于一个包含 n 条指令的序列,其需要的执行时间是什么?

- *4.41 设计一个数据结构,离线实现习题 4.40 中的操作。离线执行包含 n 条这类指令的序列,需要的执行时间为多少?
- *4.42 在算法 4.5 中,假定初始划分 $\pi = \{B_1, \dots, B_q\}$ 具有这样的属性:对于每个 $1 \leq i \leq q$,有某个 j ,使 $f^{-1}(B_i) \subseteq B_j$ 。证明:在图 4-35 的第 6 行,至多能选择一个 j 。通过这一简化,算法 4.5 的时间复杂度减少了吗?

研究性问题

- 4.43 从 4.1 节中给定的 7 个操作或者其他相关的原语中选择 n 个操作执行序列,关于可以多快地执行这样的序列,还存在许多未决的问题。如果从任意的集合中选择元素,获取与它们有关的信息的唯一方式是通过比较,那么对于任意的原语集合,不管以何种方式排列,其执行时间必定为 $O_c(n \log n)$ 。然而,当全集为整数集 $\{1, 2, \dots, n\}$ (或者对于固定的 k , 范围为 $1 \sim n^k$) 时,或者当原语集合不是充分有序时,对于需要多于 $O(n)$ 时间的东西,则鲜有评论。因此,对于图 4-36 所给出的期望或最坏情况时间,还有许多地方亟待改进。换句话说,基于整数 $1 \sim n$,对于原语的某个子集(或全部),能够给出优于 $O(n)$ 的下界吗?

文献及注释

关于散列表技术的文献有 Morris[1968]、Aho 和 Ullman[1973]以及 Knuth[1973a]。后者也包含了关于二叉查找树的信息。构建静态二叉查找树的算法 4.2 源自 Gilbert 和 Moore[1959]。对于相同问题的 $O(n^2)$ 算法可以在 Knuth[1971]中找到,也可以在其中找到习题 4.6 的解答。Hu 和 Tucker[1971]证明了用 $O(n^2)$ 时间和 $O(n)$ 空间足以构建一棵最优二叉查找树,其数据仅出现在叶子上。Knuth[1973a]证明了可在 $O(n \log n)$ 时间内实现该算法。

Reingold[1972]给出了许多基本的集合操作,如并和交的最优算法。显然, M. D. McIlroy 和 R. Morris 最先使用了算法 4.3 中的 UNION_FIND 操作。Knuth[1969]认为路径压缩是由 A. Titter 提出的。定理 4.4 来自 Hopcroft 和 Ullman[1974], 定理 4.5 则由 Tarjan[1974]提出。4.8 节中讨论的标识符等价性和替换计算的应用来自 Galler 和 Fischer[1964]。关于有穷自动机等价性的应用 3 出自 Hopcroft 和 Karp[1971]。习题 4.16 和 4.17 是关于未使用路径压缩的算法 4.3 的复杂性,分别出自于 Fischer[1972]和 Paterson[1973]。

Adel'son-Vel'skii 和 Landis[1962]提出了 AVL 树。可以在 Crane[1972]中找到习题 4.31 和 4.32 的解答。有界平衡树的概念来自 Nievergelt 和 Reingold[1973], 对于习题 4.34 ~ 4.37, 读者可以查阅该参考书。可以在 Ullman[1974]中找到使用 2-3 树的一种扩展。

习题 4.22 对 UNION-FIND 问题给出了一个早期的解决方案,这可以在 Stearns 和 Rosenkrantz[1969]中找到。习题 4.38 则出自 Aho、Hopcroft 和 Ullman[1974]。

在线和离线算法的差异可参见 Hartmanis、Lewis 和 Stearns[1965]。Rabin[1963]研究了一种重要的受限在线计算形式,称为实时计算。

划分算法取自 Hopcroft[1971], 有限自动机中状态最小化的应用也取自于此。

第5章 图 算 法

工程与科学领域中的众多问题都可以用无向或有向图表示。本章将讨论一些主要的图问题,其解决方案(在运行时间上)与图的顶点数呈多项式关系,因此也与图的边数多项式呈关系。这里将重点考虑图连通性的处理。本章包括寻找生成树、双连通分支、强连通分支以及顶点间路径的算法。第10章将研究更难的图问题。

5.1 最小代价生成树

设 $G=(V, E)$ 为连通无向图, 边依据代价函数映射到实数。在 V 中, 生成树指的是连通所有顶点的无向树。生成树的代价就是各边的代价之和。我们的目标是要为 G 找到一棵具有最小代价的生成树。一般对于具有 e 条边的图, 可以在 $O(e \log e)$ 时间内找到一棵最小代价生成树。如果较之于顶点数 e 足够大, 可在 $O(e)$ 时间内找到这样的生成树(见习题 5.3)。许多生成树算法都是基于下列两个引理。

引理 5.1 设 $G=(V, E)$ 为连通无向图, 且 $S=(V, T)$ 为 G 的一棵生成树, 则

a) 对 V 中所有的 v_1 和 v_2 , 在 S 中的 v_1 和 v_2 之间的路径是唯一的;

b) 如果将 $E-T$ 中的任何边加入 S 中, 则产生唯一回路。

证明: 如果存在多于一条路径, 则存在一个回路, 因此, (a) 是显而易见的。

因为在所添加边的两个端点之间必定已经存在一条路径, 所以 (b) 也是显然的。 \square

引理 5.2 设 $G=(V, E)$ 为连通无向图, c 为其边上的代价函数。对于 $k>1$, 设 $\{(V_1, T_1), (V_2, T_2), \dots, (V_k, T_k)\}$ 为 G 的任意生成森林。令 $T=\bigcup_{i=1}^k T_i$ 。假定 $e=(v, w)$ 为 $E-T$ 中代价最小的一条边, 满足 $v \in V_1$ 和 $w \notin V_1$, 则存在 G 的一棵生成树, 包括 $T \cup \{e\}$, 其代价与 G 中包含 T 的任何生成树一样低。

证明: 反之, 假定 $S'=(V, T')$ 为 G 的一棵生成树, 使 T' 包含 T , 但不含 e , 且 S' 的代价低于 G 中含有边集 $T \cup \{e\}$ 的任何生成树。

由引理 5.1(b), 加入 e 到 S' 中形成一个回路, 如图 5-1 所示。该回路必然包含不同于边 e 的边 $e'=(v', w')$, 使得 $v' \in V_1$ 且 $w' \in V_1$ 。由假设, $c(e) \leq c(e')$ 。

考虑图 S , 它通过向 S' 中加入边 e 并删除边 e' 而得到。由于删除边 e' 时, 破坏了唯一的环, 所以 S 无环。而且, 在 S 中 v' 与 w' 之间存在一条路径, 因而, V 中的所有顶点仍然是连通的。因此, S 是 G 的一棵生成树。因 $c(e) \leq c(e')$, S 的代价不会高于 S' 。但是, S 包含 T 和 e , 与 S' 最小相矛盾。 \square

现在给出一个为无向图 $G=(V, E)$ 找到一棵最小代价生成树的算法。该算法本质上与例 4.1 中的算法类似。该算法维护着一个由不相交的顶点集所组成的集合 VS 。VS 中的每个集合 W 皆表示一个顶点的连通集, 在用 VS 表示的生成森林中, 这些顶点组成一棵生成树。算法以代价增序, 从 E 中选择边。依次考虑每条边 (v, w) 。如果 v 和 w 已经处于 VS 中的相同集合内, 则放弃这条边。如果 v 和 w 位于不同的集合 W_1 和 W_2 (意味着 W_1 和 W_2 尚未连通) 中, 则将 W_1 和 W_2 合并到一个集合, 并将 (v, w) 加入最终生成树的

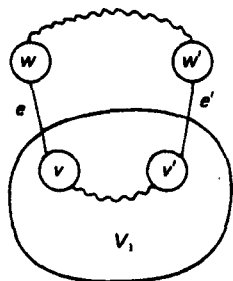


图 5-1 图 G 中的回路

边集合 T 中。可以使用 4.7 节中的不相交集合并的算法。由引理 5.2 以及对所选择边数的简单归纳可知, 对于 G , 至少有一棵最小代价生成树将包含这条边。

算法 5.1 最小代价生成树 (Kruskal 算法)。

输入: 无向图 $G = (V, E)$, 以及边上的代价函数 c 。

输出: G 的最小代价生成树 $S = (V, T)$ 。

方法: 图 5-2 给出的程序。

□

```

begin
1.   $T \leftarrow \emptyset$ ;
2.   $VS \leftarrow \emptyset$ ;
3.  construct a priority queue  $Q$  containing all edges in  $E$ ;
4.  for each vertex  $v \in V$  do add  $\{v\}$  to  $VS$ ;
5.  while  $\|VS\| > 1$  do
      begin
6.          choose  $(v, w)$ , an edge in  $Q$  of lowest cost;
7.          delete  $(v, w)$  from  $Q$ ;
8.          if  $v$  and  $w$  are in different sets  $W_1$  and  $W_2$  in  $VS$  then
              begin
9.                  replace  $W_1$  and  $W_2$  in  $VS$  by  $W_1 \cup W_2$ ;
10.                 add  $(v, w)$  to  $T$ 
              end
          end
      end
end
end

```

图 5-2 最小代价生成树算法

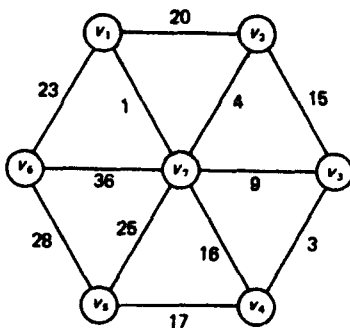


图 5-3 包含边代价的无向图

例 5.1 考虑图 5-3 的无向图。各条边按代价的递增顺序排列如下:

边	代价	边	代价
(v_1, v_7)	1	(v_4, v_8)	17
(v_8, v_4)	3	(v_1, v_2)	20
(v_2, v_7)	4	(v_1, v_8)	23
(v_8, v_7)	9	(v_5, v_7)	25
(v_8, v_3)	15	(v_8, v_6)	28
(v_4, v_7)	16	(v_6, v_7)	36

当然, 实际上并未按算法 5.1 中的步骤 3 对边排序, 而是将它们保存在一个堆、2-3 树或是其他合适的数据结构中, 一直到需要它们的时候为止。事实上, 实现优先队列, 3.4 节中的堆是

一个理想的选择。在第 6 行中, 重复查询最小的代价边是堆排序的基本操作。而且, 为了构造一棵生成树, 在第 6 行选择的边的数目通常小于 $\|E\|$ 。在这类情况下, 由于没有对 E 进行完全的排序, 从而节省了时间。

最初, 在 VS 内, 每个顶点自身位于一个集合中。最低代价的边为 (v_1, v_7) , 因此, 将该边添加到树中, 合并 VS 中的集合 $\{v_1\}$ 和 $\{v_7\}$ 。然后考虑边 (v_3, v_4) 。由于 v_3 和 v_4 位于 VS 内的不同集合中, 将 (v_3, v_4) 加入树中, 并且合并集合 $\{v_3\}$ 和 $\{v_4\}$ 。然后, 加入 (v_2, v_7) , 并合并集合 $\{v_2\}$ 与 $\{v_1, v_7\}$ 。再加入第 4 条边 (v_3, v_7) , 并合并 $\{v_1, v_2, v_7\}$ 和 $\{v_3, v_4\}$ 。

接下来考虑 (v_2, v_3) 。 v_2 和 v_3 位于同一个集合 $\{v_1, v_2, v_3, v_4, v_7\}$ 中。因此, 存在一条从 v_2 到 v_3 的路径, 包含已经位于生成树中的边, 所以没有必要加入 (v_2, v_3) 。图 5-4 中简要给出了整个的步骤序列。最终得到的无向生成树如图 5-5 所示。□

边	动作	VS 中的集合 (连通分支)
(v_1, v_7)	Add	$\{v_1, v_7\}, \{v_2\}, \{v_3\}, \{v_4\}, \{v_5\}, \{v_6\}$
(v_3, v_4)	Add	$\{v_1, v_7\}, \{v_2\}, \{v_3, v_4\}, \{v_5\}, \{v_6\}$
(v_2, v_7)	Add	$\{v_1, v_2, v_7\}, \{v_3, v_4\}, \{v_5\}, \{v_6\}$
(v_3, v_7)	Add	$\{v_1, v_2, v_3, v_4, v_7\}, \{v_5\}, \{v_6\}$
(v_2, v_3)	Reject	
(v_4, v_7)	Reject	
(v_4, v_5)	Add	$\{v_1, v_2, v_3, v_4, v_5, v_7\}, \{v_6\}$
(v_1, v_2)	Reject	
(v_1, v_6)	Add	$\{v_1, \dots, v_7\}$

图 5-4 构造生成树的步骤序列

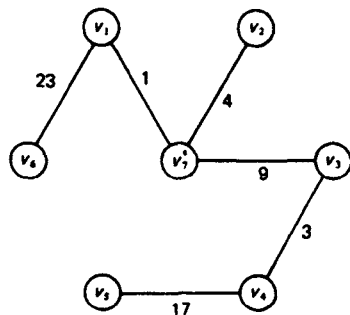


图 5-5 最小代价生成树

定理 5.1 如果图 G 是连通的, 则用算法 5.1 查找最小代价生成树。若在第 5~10 行的循环中验证 d 条边, 所消耗的时间为 $O(d \log e)$, 其中 $e = \|E\|$ 。因此, 算法 5.1 至多需要 $O(e \log e)$ 时间。

证明: 第 8 和 9 行正确地使得顶点位于相同的集合中, 当且仅当它们位于用 VS 表示的生成森林的同一棵树中。由该事实及引理 5.2, 可判定算法的正确性。

对于时间, 假定第 5~10 行循环需要 d 次迭代。如果用优先队列实现 Q , 则在 Q 中查找一条最小代价边 (第 6 行) 需要 $O(\log e)$ 步。如果使用不相交集合并的快速算法, 找到所有包含 v 和 w 的集合 W_1 和 W_2 (第 8 行), 并用它们的并取代它们 (第 9 行) 所需要的总时间至多为 $O(eG(e))$ 。显然, 循环的余下部分要求与 G 的大小不相关的常数时间。 Q 的初始化需要 $O(e)$ 时间, VS 的初始化需要 $O(n)$ 时间, 其中 n 为 V 中的顶点数。□

5.2 深度优先搜索

考虑按如下方式访问一个无向图的顶点。选定并“访问”起始顶点 v 。然后选择任意以 v 为起点的边 (v, w) , 并访问 w 。一般假定 x 为最新访问过的顶点。选择以 x 为起点, 且之前未曾出现的边 (x, y) , 继续搜索。如果 y 之前访问过, 则找到另一条以 x 为起始点的新边。如果 y 之前未访问过, 那么访问 y 并以 y 为起点进行新的搜索。在通过所有以 y 为起始的路径完成搜索之后, 搜索返回到最先到达 y 的顶点 x 。选择以 x 为起点且未曾出现的边, 继续该过程直到列表中的边都已用尽。这种访问无向图顶点的方法称为深度优先搜索, 因为该搜索尽可能以前向 (更深) 方式进行。

也可以将深度优先搜索应用到一个有向图中。如果图是有向的, 则在顶点 x 处仅选择从 x 出发的边 (x, y) 。在耗尽所有从 y 出发的边之后, 即使存在其他指向 y 且还没有搜索到的边, 也返

回到 x 。

如果将深度优先搜索应用到连通无向图，则可以看到图中每个顶点都将被访问，且每条边都将被检查。如果图不是连通的，则搜索到的将是图的连通分支。当搜索完一个连通分支之后，选择一个尚未访问过的顶点作为新的起始顶点，并开始新的搜索。

无向图 $G = (V, E)$ 的深度优先搜索将 E 中的边划分为两个集合 T 和 B 。当处于顶点 v 考虑边 (v, w) 时，如果顶点 w 之前尚未访问，则边 (v, w) 放入集合 T 中。否则，边 (v, w) 放入集合 B 中。集合 T 中的边称为树边 (tree edges)，而 B 中的边则称为回向边 (back edges)。子图 (V, T) 为一个无向森林，称为 G 的深度优先生成森林。在此情况下，森林包含一棵树 (V, T) ，称为深度优先生成树。注意，如果 G 是连通的，深度优先生成森林将是一棵树。该深度优先生成森林中的每棵树将认为是有根的，树的深度优先搜索开始于根所在的顶点。

下面给出图的深度优先搜索算法。

算法 5.2 无向图的深度优先搜索。

输入：图 $G = (V, E)$ ，由邻接表 $L[v]$ 表示，其中 $v \in V$ ；

输出：边集 E 的划分：树边集合 T 和回向边的集合 B ；

方法：在通过起始于 v 的一条边进行搜索期间，如果最先到达顶点 w ，则图 5-6 中的递归程序 $\text{SEARCH}(v)$ 将边 (v, w) 加入集合 T 。假定所有的顶点初始时都标记为“新”(new)。整个算法如下：

```

begin
6.    $T \leftarrow \emptyset$ ;
7.   for all  $v$  in  $V$  do mark  $v$  "new";
8.   while there exists a vertex  $v$  in  $V$  marked "new" do
9.      $\text{SEARCH}(v)$ 
end

```

end

所有在 E 中但不在 T 中的边都认为位于 B 中。注意，如果边 (v, w) 位于 E 中，则 w 将处于 $L[v]$ 中，且 v 将位于 $L[w]$ 中。因此，如果当前位于顶点 v ，且由于顶点 w 可能为 v 的父亲而标记为“旧”(old)，则不能够简单地将边 (v, w) 置于 B 中。□

例 5.2 为了方便，将 T 中的树边用实线表示，而将 B 中的回向边用虚线表示。并且，树（或多棵树）的根（在第 8 行选定的起始顶点）将位于顶部。从左到右，按对应的边在 SEARCH 第 4 行中添加的顺序，画出每个顶点的儿子。考虑图 5-7a，该图一个可能的树边和回向边的划分导致了图 5-7b 所示的深度优先搜索。

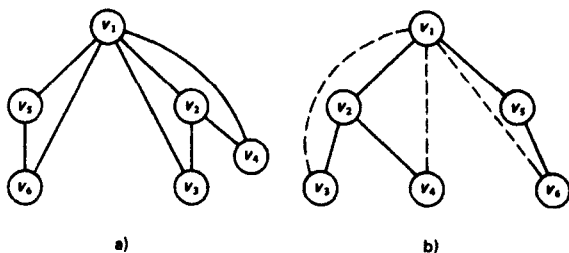


图 5-7 图及其深度优先生成树

```

procedure SEARCH(v):
begin
1.   mark  $v$  "old";
2.   for each vertex  $w$  on  $L[v]$  do
3.     if  $w$  is marked "new" then
4.       begin
5.         add  $(v, w)$  to  $T$ ;
6.          $\text{SEARCH}(w)$ 
7.       end
8.   end
end

```

图 5-6 深度优先搜索

初始时, 所有顶点都是“新”。假定在第 8 行中选定 v_1 。当执行 $\text{SEARCH}(v_1)$, 可以在第 2 行选择 $w = v_2$ 。由于 v_2 标记为“新”, 将 (v_1, v_2) 加入 T , 调用 $\text{SEARCH}(v_2)$ 。 $\text{SEARCH}(v_2)$ 可以从 $L[v_2]$ 中选择 v_1 , 但是 v_1 已经标记为“旧”。假定选定 $w = v_3$, 因为 v_3 是“新”, 将 (v_2, v_3) 添加到 T , 调用 $\text{SEARCH}(v_3)$ 。与 v_3 邻接的每个顶点现在都是“旧”, 因此返回到 $\text{SEARCH}(v_2)$ 。

然后, 继续 $\text{SEARCH}(v_2)$, 找到边 (v_2, v_4) , 将其加入 T , 并调用 $\text{SEARCH}(v_4)$ 。注意, 在图 5-7b 中, v_4 位于之前发现的 v_2 的儿子 v_3 的右边。没有“新”顶点与 v_4 邻接, 因此返回 $\text{SEARCH}(v_2)$ 。此时, 设找到“新”顶点同 v_2 邻接, 因此返回 $\text{SEARCH}(v_1)$ 。继续 $\text{SEARCH}(v_1)$, 找到 v_5 , 且由 $\text{SEARCH}(v_5)$ 找到 v_6 。那么, 所有顶点将位于树中并标记为“旧”, 因此算法终止。如果图不是连通的, 第 8~9 行中的循环将继续, 对于每个分支重复一次。□

定理 5.2 对于有 n 个顶点 e 条边的图, 算法 5.2 需要 $O(\text{MAX}(n, e))$ 步。

证明: 如果生成一个顶点列表并扫描一次, 则第 7 行以及第 8 行中对“新”顶点的搜索需要 $O(n)$ 步。除了对本身的递归调用, 在 $\text{SEARCH}(v)$ 中花费的时间同 v 所邻接的顶点数成比例。由于在 $\text{SEARCH}(v)$ 第一次调用的时候, v 标记为“旧”, 所以, 对于给定的 v , $\text{SEARCH}(v)$ 仅调用一次。因此, SEARCH 所花费的总时间为 $O(\text{MAX}(n, e))$, 从而定理得证。□

深度优先搜索的部分能力包含在下列引理中, 假定无向图 G 的每条边不是深度优先森林中的一条边, 就是在深度优先生成森林中的某棵树中连接祖先及其后代。因此, 无论是树边还是回向边, G 中所有的边连接两个顶点, 使得在生成森林中一个顶点是另一个顶点的祖先。

引理 5.3 如果 (v, w) 为一条回向边, 则在生成森林中, v 为 w 的祖先或 w 为 v 的祖先。

证明: 不失一般性, 假定在访问 w 之前访问 v , 在此意义上, $\text{SEARCH}(v)$ 在 $\text{SEARCH}(w)$ 之前调用。因此, 当到达 v 的时候, w 仍标记为“新”。在生成森林中, 通过 $\text{SEARCH}(v)$ 访问的所有“新”顶点将成为 v 的后代。但是, 由于 w 处于列表 $L[v]$ 中, $\text{SEARCH}(v)$ 不可能终止直到到达 w 。□

存在一个自然序, 将深度优先搜索施加到生成森林的顶点。换句话说, 如果在算法 5.2 的第 6 和 7 行之间将 COUNT 初始化为 1, 并在 SEARCH 程序的开始处插入如下语句, 则可以按其访问的顺序标记顶点:

$\text{DFNUMBER}[v] \leftarrow \text{COUNT};$

$\text{COUNT} \leftarrow \text{COUNT} + 1;$

然后, 森林中的顶点将标记为 1, 2, ..., 一直到森林中的顶点数。

对于 n 个顶点的图, 可以在 $O(n)$ 时间内完全地安排这些标记。在最终得到的生成森林中, 标记的顺序同每棵树的前序遍历相一致。而后将假定, 所有深度优先生成森林也是按此方式标记。通常, 即使其本身就是顶点的名字, 也将处理这些顶点标记。因此, 这使某些说法有意义, 例如, “ $v < w$ ”, 其中 v 和 w 为顶点。

例 5.3 图 5-7a 中顶点的深度优先序为 $v_1, v_2, v_3, v_4, v_5, v_6$, 通过对每个顶点追踪 SEARCH 初始化的顺序来保证这种优先序, 或通过先序方式遍历图 5-7b 中的树来确定。□

特别注意, 如果 v 是 w 的合适的祖先, 则有 $v < w$ 。另外, 在一棵树中如果 v 在 w 的左边, 则 $v < w$ 。

5.3 双连通性

现在考虑应用深度优先搜索来确定无向图的双连通分支。设 $G = (V, E)$ 为连通无向图。如果存在顶点 v 和 w , 使 v 和 w 之间的每条路径都包含相异于 v, w 的顶点 a , 则顶点 a 称为 G 的关节点 (articulation point)。换句话说, 若删除 a 将把 G 分解为两个或更多的部分, 则 a 为 G 的一个关节点。如果对于相异的 3 个顶点 v, w 和 a , 在 v 和 w 之间存在一条不包含 a 的路径, 则称图 G

是双连通的。因此, 一个无向连通图为双连通的, 当且仅当其无关节点。

如果两条边 $e_1 = e_2$ 或者存在含有 e_1 和 e_2 的环, 则称边 e_1 和 e_2 关联, 此时能够在 G 的边集上定义一个自然关系。容易证明, 该关系为一个等价关系^①, 它将 G 的边划分为等价类 E_1, E_2, \dots, E_k , 满足: 两条不同的边处于同一类中, 当且仅当它们处在一个公共环中。对于 $1 \leq i \leq k$, 设 V_i 为相应于 E_i 中边的顶点集, 称每个图 $G_i = (V_i, E_i)$ 为 G 的双连通分支。

例 5.4 考虑图 5-8a 中的无向图。例如, 由于 v_1 和 v_7 之间的每条路径都经由 v_4 , 所以顶点 v_4 为一个关节点。公共环上的等价类为

$$\begin{aligned} & \{(v_1, v_2), (v_1, v_3), (v_2, v_3)\}, \\ & \{(v_2, v_4), (v_2, v_5), (v_4, v_5)\}, \\ & \{(v_4, v_6)\}, \\ & \{(v_6, v_7), (v_6, v_8), (v_6, v_9), \\ & (v_7, v_8), (v_8, v_9)\}. \end{aligned}$$

这些集合引出了图 5-8b 所示的双连通分支。唯一的非直观处是边 (v_4, v_6) 本身是等价类 (不存在包含这条边的环), 其所在的“双连通分支”仅包含顶点 v_4 和 v_6 。

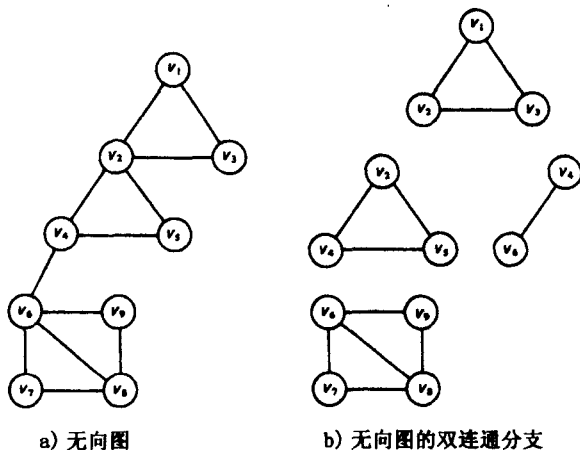


图 5-8

下述引理给出一些与双连通性相关的有用信息。

引理 5.4 对于 $1 \leq i \leq k$, 设 $G_i = (V_i, E_i)$ 为连通无向图 $G = (V, E)$ 的双连通分支, 那么,

1. 对于每个 $i (1 \leq i \leq k)$, G_i 是双连通的。
2. 对于所有的 $i \neq j$, $V_i \cap V_j$ 至多包含一个顶点。
3. a 为 G 的一个关节点, 当且仅当对于某 $i \neq j$, 有 $a \in V_i \cap V_j$ 。

证明:

1. 假定 V_i 中存在 3 个相异的顶点 v, w 和 a , 使 G_i 中所有 v 和 w 之间的路径都经由 a , 那么, 肯定 (v, w) 不是 E_i 中的一条边。因此, E_i 中存在不同的边 (v, v') 和 (w, w') , 并且在 G_i 中存在这些边的环。由双连通分支的定义, 该环的所有边和顶点都分别位于 E_i 和 V_i 中。因此, 在 G_i 中 v 和 w 之间存在两条路径, 仅有一条包含 a , 矛盾。

2. 假定两个相异顶点 v 和 w 处于 $V_i \cap V_j$ 中。那么, 在 G_i 中存在一个包含 v 和 w 的环 C_1 , 在 G_j 中也存在一个包含 v 和 w 的环 C_2 。由于 E_i 和 E_j 不相交, 所以, 在 C_1 和 C_2 中的边集是不相交的。然而, 可以通过使用 C_1 和 C_2 中的边构造一个包含 v 和 w 的环, 意味着在 E_i 中至少有一条边等价于 E_j 中的一条边。因此, 按照假定的条件, 将使得 E_i 和 E_j 不是等价类。

3. 假定顶点 a 是 G 的关节点, 则存在两个顶点 v 和 w , 满足 v, w 和 a 相异, 且 v 和 w 之间的每条路径都包含 a 。由于 G 是连通的, 至少存在一条这样的路径。设 (x, a) 和 (y, a) 为 v 和 w 之间一条与 a 关联的路径上的两条边。如果存在一个环含有这两条边, 那么 v 和 w 之间存在一条不包含 a 的路径。因此, (x, a) 和 (y, a) 处于不同的双连通分支中, 且 a 为它们顶点集合的交。

相反, 如果 $a \in V_i \cap V_j$, 则在 E_i 和 E_j 中分别存在边 (x, a) 和 (y, a) 。由于这两条边都没有出

① 如果 R 是自反的 (对于所有的 $a \in S, aRa$)、对称的 (对于所有的 $a, b \in S, aRb$ 蕴涵着 bRa) 和传递的 (aRb 和 bRc 蕴涵着 aRc), 则 R 为集合 S 的等价关系。容易证明, S 的等价关系将 S 划分为不相交的等价类。(子集 $[a] = \{b \mid bRa\}$ 称为等价类。)

现在任意一个环, 满足从 x 到 y 的每条路径都包含 a 。因此, a 是一个关节点。□

在寻找无向图的双连通分支时, 深度优先搜索特别有用。原因之一在于, 由引理 5.3, 不存在“交叉边”。也就是说, 在生成森林中, 如果顶点 v 既不是 w 的祖先, 也不是其后代, 那么可能不存在从 v 到 w 的边。

如果顶点 a 为关节点, 那么删除 a 及与 a 相关的所有边后, 图将分解为两个或更多的部分。其中一部分由 a 的儿子 s 及其在深度优先生成树中的所有后代组成。因此, 在深度优先生成树中, a 必定有一个儿子 s , 在 s 的后代和 a 的某个祖先之间不存在回向边。反之, 除了生成树的根, 如果不存在从 a 的任意儿子的后代到 a 的某个祖先的回向边, 则没有交叉边意味着顶点 a 为一个关节点。深度优先生成树的根是一个关节点, 当且仅当它有两个或多个儿子。

例 5.5 图 5-8a 中的图对应的深度优先生成树如图 5-9 所示。关节点为 v_2 、 v_4 和 v_6 。顶点 v_2 有儿子 v_4 , 并且从 v_4 的后代到 v_2 的祖先之间不存在回向边。同样, v_4 有儿子 v_6 , v_6 有儿子 v_8 , 都具有相似的属性。□

先前的思想体现在下列引理中。

引理 5.5 设 $G=(V, E)$ 为连通无向图, 并设 $S=(V, T)$ 为 G 的深度优先生成树。顶点 a 为 G 的一个关节点, 当且仅当

1. a 为根, 并且 a 有超过一个儿子, 或者
2. a 不是根, 并且对于 a 的某个儿子 s , 在 s 的任意后代 (包括其本身) 和 a 的某个祖先之间不存在回向边。

证明: 容易证明, 根为关节点当且仅当其有多于一个儿子。此部分留作练习。

假定条件 2 是成立的。设 f 为 a 的父亲。由引理 5.3, 每条回向边都是从一个顶点到其祖先。因此, s 的后代 v 的任意回向边指向 v 的祖先。由引理的假设, 回向边不能到达 a 的某个祖先。因此, 要么指向 a , 要么指向 s 的后代。因此, 从 s 到 f 的每条路径都包含 a , 意味着 a 为一个关节点。

为反向证明, 假定 a 为一个关节点但不是根。设 x 和 y 为相异于 a 的顶点, 满足在 G 中每条 x 和 y 之间的边都包含 a 。 x 和 y 中至少有一个假设为 x , 是 S 中 a 的某个后代, 否则, 在 G 中 x 和 y 之间存在一条路径使用 T 中的边但不包含 a 。设 s 为 a 的儿子, 使得 x 为 s 的后代 (可能 $x=s$)。或者在 s 的后代 v 和 a 的某个祖先 w 之间不存在回向边, 这直接使得条件 2 的情形是正确的, 或者存在这样的一条边 (v, w) 。在后一种情况下, 必须考虑两种情形。

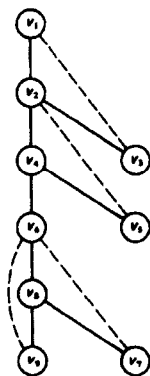


图 5-9 深度优先生成树

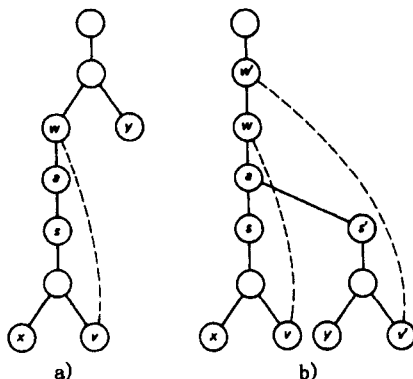


图 5-10 路径反例

情形 1 假定 y 不是 a 的后代。那么, 存在一条从 x 到 v , 再到 w , 再到 y 的路径, 避开了 a , 矛盾。见图 5-10a。

情形 2 假定 y 为 a 的后代。肯定 y 不是 s 的后代, 否则存在从 x 到 y 的一条路径, 避开 a 。

设 s' 为 a 的儿子, 使得 y 为 s' 的后代。在 s' 的后代 v' 和 a 的某个祖先 w' 之间不存在回向边, 在此情形, 条件 2 直接成立, 或者存在这样的一条边 (v', w') 。在后一种情况下, 存在从 x 到 v , 再到 w , 再到 w' , 再到 v' , 再到 y 的一条路径, 避开 a , 矛盾。见图 5-10b。至此得到结论, 条件 2 是正确的。□

对连通无向图 $G = (V, E)$ 施行深度优先搜索, 设 T 和 B 分别为由此产生的树边和回向边的集合。假定 V 中的顶点以它们的深度优先搜索数命名。对于 V 中的每个 v , 定义

$$\text{LOW}[v] = \text{MIN}(\{v\} \cup \{w \mid \text{存在回向边 } (x, w) \in B, \text{ 使得} \\ \text{在深度优先生成森林 } (V, T) \text{ 中,} \\ x \text{ 为 } v \text{ 的后代, 且 } w \text{ 为 } v \text{ 的祖先}\}) \quad (5-1)$$

先序编号意味着, 如果 x 为 v 的后代且 (x, w) 为回向边使得 $w < v$, 那么 w 为 v 的某个祖先。因此, 由引理 5.5, 如果顶点 v 不为根, 那么 v 为关节点当且仅当 v 有儿子 s 使得 $\text{LOW}[s] \geq v$ 。

如果根据经由回向边邻接到 v 的顶点和 v 的儿子的 LOW 值, 重写式 (5-1) 以表达 $\text{LOW}[v]$, 则可以在程序 SEARCH 中嵌入一个计算步确定每个顶点的 LOW 值。特别地, 要计算 $\text{LOW}[v]$, 可以通过确定顶点 w 的最小值, 使得

1. $w = v$, 或者
2. $w = \text{LOW}[s]$ 且 s 为 v 的儿子, 或者
3. (v, w) 为 B 中的回向边。

只要穷举与 v 邻接的列表 $L[v]$ 的值, 就可以确定 w 的最小值。因此式 (5-1) 等价于

$$\text{LOW}[v] = \text{MIN}(\{v\} \cup \{\text{LOW}[s] \mid s \text{ 为 } v \text{ 的儿子}\} \cup \{w \mid (v, w) \in B\}) \quad (5-2)$$

至此, 已经按优先访问对顶点进行重命名并计算 LOW , 并将其加入到图 5-11 所示 SEARCH 的修改版中。在第 4 行, 将 $\text{LOW}[v]$ 初始化为其最大可能值。如果在深度优先生成森林中, 顶点 v 有儿子 w , 要是 $\text{LOW}[w]$ 小于 $\text{LOW}[v]$ 的当前值, 则在第 11 行调整 $\text{LOW}[v]$ 。如果通过一条回向边将顶点 v 连接到顶点 w , 要是顶点 w 的深度优先数小于 $\text{LOW}[v]$ 的当前值, 则在第 13 行使 $\text{LOW}[v]$ 为 $\text{DFNUMBER}[w]$ 。由于在深度优先生成树中, w 是 v 的父亲, 所以在第 12 行中的判定验证了 (v, w) 并非真正的回向边。因此, 图 5-11 实现了方程 (5-2)。

```

procedure SEARCHB(v):
begin
1.  mark v "old";
2.  DFNUMBER[v] ← COUNT;
3.  COUNT ← COUNT + 1;
4.  LOW[v] ← DFNUMBER[v];
5.  for each vertex w on L[v] do
6.      if w is marked "new" then
7.          begin
8.              add (v, w) to T;
9.              FATHER[w] ← v;
10.             SEARCHB(w);
11.             if LOW[w] ≥ DFNUMBER[v] then a biconnected
                component has been found;
12.             LOW[v] ← MIN(LOW[v], LOW[w])
13.         end
14.     else if w is not FATHER[v] then
15.         LOW[v] ← MIN(LOW[v], DFNUMBER[w])
16. end

```

图 5-11 包含计算 LOW 值的深度优先搜索

由于已经为每一个顶点 v 找到了 $\text{LOW}[v]$, 可以很容易地识别出关节点。下面首先给出完整的算法, 然后证明其正确性, 再说明其时间开销为 $O(e)$ 。

算法 5.3 找出双连通分支。

输入：连通无向图 $G=(V, E)$ 。

输出： G 的每个双连通分支的边列表。

方法：

1. 初始时，将 T 设置为 \emptyset ，COUNT 设为 1。并将 V 中的每个顶点标记为“新”。然后选择 V 中的任意顶点 v_0 ，调用 SEARCHB(v_0) (图 5-11)，构建一棵深度优先生成树 $S=(V, T)$ ，为 V 中的每个 v 计算 $LOW(v)$ 。

2. 当在 SEARCHB 的第 5 行碰到顶点 w 时，要是边 (v, w) 不在边的下推存储 STACK 中，则将其放入 STACK 中[⊖]。在第 10 行找到 (v, w) 使得 w 为 v 的儿子且 $LOW[w] \geq v$ 后，从 STACK 中弹出所有的边直到含有 (v, w) 。这些边组成了 G 的双连通分支。□

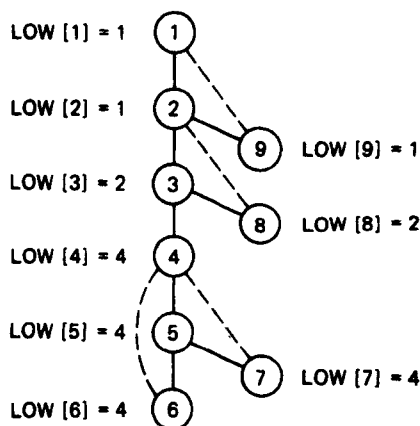


图 5-12 考虑 LOW 值的图 5-9 的生成树

例 5.6 重新生成图 5-9 的深度优先生成树如图 5-12 所示，顶点由 DFNUMBER 重命名，且显示 LOW 值。例如，由于存在回向边 $(6, 4)$ ，SEARCHB(6) 确定 $LOW[6]=4$ 。然后，由于 4 小于 $LOW[5]$ 的初始值 5，通过调用 SEARCHB(6)，SEARCHB(5) 设置 $LOW[5]=4$ 。

在完成 SEARCHB(5) 之后，发现 (第 10 行) $LOW[5]=4$ 。因此，4 为关节点。此处，下推式存储包含边 (从底至顶)：

$(1, 2), (2, 3), (3, 4), (4, 5), (5, 6), (6, 4), (5, 7), (7, 4)$

因此，往下弹出边以包含 $(4, 5)$ ，也就是说，输出边 $(7, 4)$ 、 $(5, 7)$ 、 $(6, 4)$ 、 $(5, 6)$ 和 $(4, 5)$ 为找到的第一个双连通分支的边。

观察 SEARCHB(2) 完成的工作，可以发现 $LOW[2]=1$ ，并且即使 1 不是关节点也会清空下推式存储。这确保包含根的双连通分支也可弹出。□

定理 5.3 算法 5.3 正确地找到 G 的双连通分支，如果 G 有 e 条边则需要 $O(e)$ 时间。

证明：步骤 1 需要 $O(e)$ 时间的证明是对 SEARCH (定理 5.2) 观察的一个简单扩展。步骤 2 对每条边验证一次，将其放入一个下推式存储，而后弹出它。因此步骤 2 是 $O(e)$ 的。

对于算法的正确性，引理 5.5 确保正确地找到关节点。即使根不是一个关节点，但为了弹出包含根的双连通分支，它也被认为是关节点。

必须证明，如果 $LOW[w] \geq v$ ，那么当完成 SEARCHB(w) 时，STACK 中 (v, w) 以上的边将

⊖ 注意，如果 (v, w) 为一条边， v 处于 $L[w]$ 中， w 处于 $L[v]$ 中，那么 (v, w) 要碰到两次，一次是顶点 v 访问，一次是顶点 w 访问。通过测试，如果 $v < w$ 且 w 为“旧”，或者 $v > w$ 且 $w = FATHER[v]$ ，可以验证是否 (v, w) 已经位于 STACK 中。

恰好是在含有 (v, w) 的双连通分支中的那些边。这可以通过基于 G 的双连通分支数 b 的归纳来做到。对于归纳的基 $b=1$, 显然成立, 因为在这种情况下, v 为树的根, (v, w) 为以 v 为起点的唯一树边, 并且当 $\text{SEARCHB}(w)$ 完成时, G 中所有的边都在 STACK 中。

现在, 假定归纳假设对所有包含 b 个双连通分支的图成立, 并设 G 为具有 $b+1$ 个双连通分支的图。对于树边 (v, w) , 设 $\text{SEARCHB}(w)$ 为 SEARCHB 的第一次调用, 以 $\text{LOW}(w) \geq v$ 结束。由于没有从 STACK 中删除任何边, 在 STACK 中位于 (v, w) 之上的边的集合为所有以 w 的后代作为起点的边的集合。容易证明, 这些边恰好是含有 (v, w) 的双连通分支中的那些边。当从 STACK 中删除这些边的时候, 算法的执行如在图 G' 时一样, 图 G' 由 G 通过删除含有边 (v, w) 的双连通分支而得到。现在, 可适用归纳步, 因为 G' 有 b 个双连通分支。□

5.4 有向图的深度优先搜索

如果将与 v 相“邻接”的顶点列表 $L[v]$ 定义为 $\{w \mid (v, w) \text{ 为一条边} \}$, 即 $L[v]$ 为以 v 为尾的边头顶点的列表, 那么能够用算法 5.2 为有向图 $G=(V, E)$ 寻找有向生成森林。

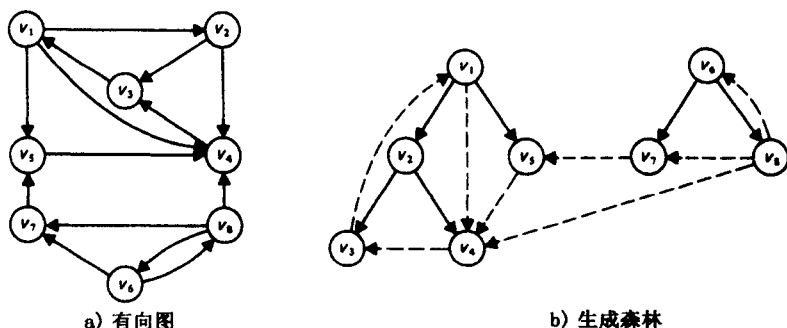


图 5-13 有向图的深度优先搜索

例 5.7 图 5-13a 给出一个有向图, 在图 5-13b 是它的深度优先生成森林。同先前一样, 树边用实线表示, 其他边用虚线表示。

为构造生成森林, 从 v_1 开始。 $\text{SEARCH}(v_1)$ 调用 $\text{SEARCH}(v_2)$, $\text{SEARCH}(v_3)$ 调用 $\text{SEARCH}(v_3)$ 。当不能继续添加到树时, 较后的过程终止。因为以 v_3 为尾端的唯一一边指向 v_1 , 而 v_1 已经标记为“旧”。因此, 返回到 $\text{SEARCH}(v_2)$, 添加 v_4 为 v_2 的第二个儿子。 $\text{SEARCH}(v_4)$ 终止, 因为 v_3 已经是“旧的”, 所以对森林不做添加。然后, $\text{SEARCH}(v_2)$ 终止, 因为除了 v_2 外, 已经考虑了所有的边。再回到 v_1 , 调用 $\text{SEARCH}(v_5)$ 。后者也以对树无添加操作而终止, 相似地, $\text{SEARCH}(v_1)$ 也不再进行任何添加。

现在, 选择 v_6 作为新深度优先生成树的根。其结构与前述相似, 留给读者证明。注意, 选择访问顶点的顺序是 v_1, v_2, \dots, v_8 。因此, 顶点 v_i 的深度优先编号为 i , $1 \leq i \leq 8$ 。□

注意, 在有向图的深度优先搜索中, 除了树边还存在 3 种类型的边。如图 5-13b 中诸如 (v_3, v_1) 的回向边, (v_4, v_3) 的从右到左交叉边以及 (v_1, v_4) 的前向边。然而, 没有任何边以较低深度优先编号的顶点为起点, 而指向深度优先编号更高的顶点, 除非后者是前者的后代。这并非偶然。

这种解释同在无向图中不存在交叉边的理由是相似的。假定 (v, w) 是一条边, v 在 w 之前被访问 (也就是 $v < w$)。在 $\text{SEARCH}(v)$ 开始到终止的那段时间, 给每个顶点安排一个编号, 成为 v 的后代。但是, 必须在边 (v, w) 出现时给 w 安排一个编号, 除非已经给 w 安排了编号。如果在边 (v, w) 出现时给 w 安排一个编号, 则 (v, w) 成为一条树边。否则, (v, w) 为一条前向边。因此可能不存在使 $v < w$ 的交叉边 (v, w) 。

在有向图 G 中, G 的深度优先搜索将边划分为 4 类:

1. 树边, 在搜索期间导向新顶点的边;
2. 前向边, 源自祖先, 并指向某个后代, 但不是树边;
3. 回向边, 源自后代而指向祖先(可能由一个顶点到其本身);
4. 交叉边, 在彼此不互为祖先或者后代的顶点之间的边。

下面的引理说明了交叉边的关键属性。

引理 5.6 如果 (v, w) 为交叉边, 则 $v > w$, 也就是说, 交叉边从右到左。

证明: 证明类似于引理 5.3, 留给读者作为练习。 \square

5.5 强连通性

对有向图进行深度优先搜索可使一些有效算法的设计成为可能, 作为例子, 考虑判定一个有向图是否强连通的问题, 即从每个顶点到任意其他顶点是否存在一条路径。

定义 设 $G = (V, E)$ 为有向图。可以将 V 划分为等价类 $V_i, 1 \leq i \leq r$, 使顶点 v 和 w 等价, 当且仅当存在一条从 v 到 w 的路径以及一条从 w 到 v 的路径。设 $E_i (1 \leq i \leq r)$ 为连接 V_i 中顶点对的边的集合。图 $G_i = (V_i, E_i)$ 称为 G 的强连通分支。即使 G 中的每个顶点位于某个 V_i 中, G 也可以有不在 E_i 中的边。一个图认为是强连通的, 如果它仅有一个强连通分支。

现在利用深度优先搜索找出图的强连通分支。首先证明每个强连通分支的顶点都是由深度优先搜索而确定的生成森林的一个连通子图。该连通子图为一棵树, 且树的根称为强连通分支的根。然而, 并非深度优先生成森林中的每棵树都必须表示一个强连通分支。

引理 5.7 设 $G_i = (V_i, E_i)$ 为有向图 G 的强连通分支。设 $S = (V, T)$ 为 G 的深度优先生成森林。那么, G_i 的顶点连同 E_i 和 T 公共的边一起形成一棵树。

证明: 设 v 和 w 为 V_i 中的顶点(假定用深度优先编号命名顶点)。不失一般性, 假定 $v < w$ 。由于 v 和 w 位于相同的强连通分支中, 在 G_i 中存在一条从 v 到 w 的路径 P 。设 x 为 P 中最小编号的顶点(可能是 v 本身)。一旦路径 P 到达 x 的一个后代, 就不可能离开 x 后代所在的子树, 因为唯一能够离开子树的边是指向编号小于 x 的顶点的交叉边和回向边。(因为从 x 开始, x 的后代连续编号。离开 x 后代的子树的交叉边或回向边必须指向编号小于 x 的顶点。)因此, w 为 x 的后代。由于顶点按先序编号, 所有编号在 x 和 w 之间的顶点也是 x 的后代。由于 $x \leq v < w$, v 为 x 的后代。

至此, 已经证明 G_i 中的任意两个顶点在 G_i 中有相同的祖先。设 r 为 G_i 中顶点的最小编号公共祖先。如果 v 在 G_i 中, 那么在生成树中从 r 到 v 路径上的任意顶点也位于 G_i 中。 \square

在有向图 G 的深度优先搜索期间, 通过按根最后被碰到的顺序找出各分支的根, 可以找到 G 的强连通分支。设 r_1, r_2, \dots, r_k 为顶点深度优先搜索结束时按序排列的根(也就是 r_i 的搜索在 r_{i+1} 之前结束)。然后, 对于每个 $i < j$, r_i 或者在 r_j 的左边, 或者 r_i 为深度优先生成森林中 r_j 的后代。

设 G_i 为强连通分支, 其根为 r_i , 其中 $1 \leq i \leq k$ 。那么, G_i 包含 r_i 的所有后代, 由于对于 $j > i$, r_j 不可能为 r_i 的后代。相似地, 可以证明下述引理。

引理 5.8 对于每个 $i (1 \leq i \leq k)$, G_i 所包含为 r_i 的后代的顶点, 但不在 G_1, G_2, \dots, G_{i-1} 中。

证明: 对于 $j > i$, 根 r_j 不可能为 r_i 的后代, 因为 $\text{SEARCH}(r_j)$ 的调用在 $\text{SEARCH}(r_i)$ 之后终止。 \square

为有助于找到根, 定义名为 LOWLINK 的函数:

$\text{LOWLINK}[v] = \text{MIN}(\{v\} \cup \{w \mid \text{存在一条从 } v \text{ 的后代到 } w \text{ 的交叉边或回向边, 并且含有 } w \text{ 的}\})$

强连通分支的根为 v 的祖先)

(5-3)

图 5-14 给出了一条从 v 的后代到 w 的交叉边, 其中, 含有 w 的强连通分支的根 r 为 v 的祖先。

很快将看到执行深度优先搜索时, 如何计算 LOWLINK。

首先, 根据 LOWLINK 提供强连通分支根的特点。

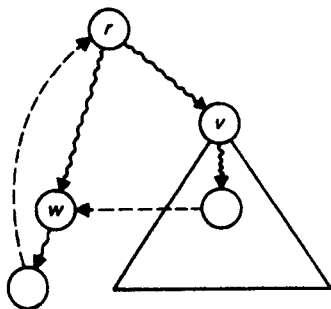
引理 5.9 设 G 为有向图。顶点 v 为 G 的强连通分支的根, 当且仅当 $\text{LOWLINK}[v] = v$ 。

证明:

“仅当”。假定 v 为 G 的强连通分支的根。由 LOWLINK 的定义, $\text{LOWLINK}[v] \leq v$ 。假设有 $\text{LOWLINK}[v] < v$ 。那么, 存在顶点 w 和 r , 使得

1. 通过一条源自 v 的后代的交叉边或回向边可以到达 w ,
2. r 为含有 w 的强连通分支的根,
3. r 为 v 的祖先, 并且
4. $w < v$ 。

图 5-14 满足 LOWLINK 条件的交叉边



由条件 2, r 为 w 的一个祖先, 因此 $r \leq w$ 。因此, 由条件 4, $r < v$, 根据条件 3, 这意味着 r 为 v 的某个祖先。但是, r 和 v 必须位于相同的强连通分支中, 因为在 G 中存在一条从 r 到 v 的路径和一条从 v 到 r 经由 w 的路径。因此, v 不是强连通分支的根, 矛盾。所以, 可以得出结论 $\text{LOWLINK}[v] = v$ 。

“当”。设 $\text{LOWLINK}[v] = v$ 。如果 v 不是含有 v 的强连通分支的根, 那么 v 的某个后代 r 为根。因此, 存在一条从 v 到 r 的路径 P 。考虑 P 的第一条边, 从 v 的后代到顶点 w , 其中 w 不是 v 的后代。这条边要么是指向 v 的祖先的一条回向边, 要么是指向编号小于 v 的顶点的一条交叉边。在这两种情况下, $w < v$ 。

还需要证明, r 和 w 位于 G 的相同强连通分支中。由于 r 是 v 的祖先, 必定存在从 r 到 v 的一条路径。路径 P 从 v 到 w , 再到 r 。因此, r 和 w 处在相同的强连通分支中。因此, $\text{LOWLINK}[v] \leq w < v$, 矛盾。□

在深度优先搜索期间, 很容易计算 LOWLINK 的值。也可以按下列方式很容易地找到强连通分支。按照搜索期间访问到的顺序, 将 G 的顶点放到栈中。无论何时找到一个根, 相关强连通分支的所有顶点都位于栈的顶端, 并且被弹出。这种策略依据引理 5.8 和深度优先编号的属性“实施”。

第一次到达顶点 v 的时候, 将 $\text{LOWLINK}[v]$ 设置为等于 v 。如果出现一条回向边或交叉边 (v, w) , 且 w 同 v 的某个祖先位于相同的强连通分支当中, 那么将 $\text{LOWLINK}[v]$ 设置为当前值和 w 中的最小值。如果出现树边 (v, w) , 则以 w 为根的子树递归出现, 计算 $\text{LOWLINK}[w]$ 。在返回 v 之后, 将 $\text{LOWLINK}[v]$ 设置为其当前值和 $\text{LOWLINK}[w]$ 中的最小值。

通过检测是否 w 仍处于顶点栈中来完成检测, 以确定是否 w 同 v 的某个后代位于相同分支中。利用一个数组表明是否顶点位于栈中来实现这种检测。注意, 由引理 5.8, 如果 w 仍处于栈中, 那么含有 w 的强连通分支的根为 v 的祖先。

现在修改程序 SEARCH, 计算 LOWLINK。对于顶点适用一个下推列表 STACK。程序如图 5-15 所示。

```

procedure SEARCHC( $v$ ):
begin
1.   mark  $v$  "old";
2.   DFNUMBER[ $v$ ]  $\leftarrow$  COUNT;
3.   COUNT  $\leftarrow$  COUNT + 1;
4.   LOWLINK[ $v$ ]  $\leftarrow$  DFNUMBER[ $v$ ];
5.   push  $v$  on STACK;
6.   for each vertex  $w$  on  $L[v]$  do
7.       if  $w$  is marked "new" then
8.           begin
9.               SEARCHC( $w$ );
               LOWLINK[ $v$ ]  $\leftarrow$  MIN(LOWLINK[ $v$ ],
               LOWLINK[ $w$ ])
           end
10.      else if DFNUMBER[ $w$ ] < DFNUMBER[ $v$ ] and  $w$  is on
               STACK then
11.          LOWLINK[ $v$ ]  $\leftarrow$  MIN(DFNUMBER[ $w$ ],
               LOWLINK[ $v$ ]);
12.      if LOWLINK[ $v$ ] = DFNUMBER[ $v$ ] then
13.          begin
               repeat
14.                 begin
                     pop  $x$  from top of STACK;
                     print  $x$ 
                 end
               until  $x = v$ ;
15.          print "end of strongly connected component"
16.      end
end

```

图 5-15 计算 LOWLINK 的程序

LOWLINK 的计算出现在第 4、9 和 11 行。在第 4 行，初始化 LOWLINK[v] 为顶点 v 的深度优先编号。在第 9 行，如果对于某个儿子 w ，LOWLINK[w] 小于 LOWLINK[v] 的当前值，那么设置 LOWLINK[v] 为 LOWLINK[w]。在第 10 行，确定 (v, w) 是回向边还是交叉边，并检测是否已经找到含有 w 的强连通分支。如果没有，那么含有 w 的强连通分支的根为 v 的祖先。必然地，在第 11 行，如果已经不存在更小的值，则将 LOWLINK[v] 设置为 w 的深度优先编号。

下面给出为有向图找到强连通分支的完整算法。

算法 5.4 有向图的强连通分支。

输入：有向图 $G = (V, E)$ 。

输出： G 的强连通分支列表。

方法：

```

begin
    COUNT  $\leftarrow$  1;
    for all  $v$  in  $V$  do mark  $v$  "new";
    initialize STACK to empty;
    while there exists a vertex  $v$  marked "new" do SEARCHC( $v$ )
end

```

□

例 5.8 考虑图 5-13 的深度优先生成森林，计算 LOWLINK 后又得到的结果如图 5-16 所示。第一个 SEARCHC 调用终止在 SEARCHC(3)。当检测边 $(3, 1)$ 时，将 LOWLINK[3] 设置为 $\text{MIN}(1, 3) = 1$ 。当返回 SEARCHC(2) 的时候，将 LOWLINK[2] 设置为 $\text{MIN}(2, \text{LOWLINK}[3]) = 1$ 。然后，调用 SEARCHC(4)，考虑边 $(4, 3)$ 。由于 $3 < 4$ ，并且 3 仍在 STACK 中，将 LOWLINK[4] 设置为 $\text{MIN}(3, 4) = 3$ 。

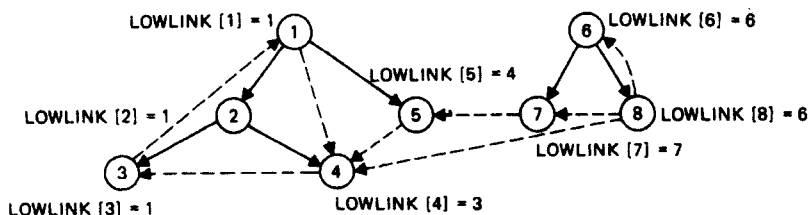


图 5-16 计算 LOWLINK 后的生成森林

然后返回 SEARCHC(2)，并将 LOWLINK[2] 设置为 LOWLINK[4] 和 LOWLINK[2] 当前值（也就是 1）中的较小值。由于后者更小，LOWLINK[2] 不改变。返回到 SEARCHC(1)，设置 LOWLINK[1] 为 $\min(1, \text{LOWLINK}[2]) = 1$ 。如果考虑边 (1, 4)，由于 (1, 4) 为一条前向边且第 10 行 SEARCHC 的条件由于 $4 > 1$ 而不满足，那么将不会有什么行为发生。

接着，调用 SEARCHC(5)，由于 $4 < 5$ 且 4 在 STACK 中，交叉边 (5, 4) 使 LOWLINK[5] 设置为 4。当再次返回 SEARCHC(1) 的时候，设置 LOWLINK[1] 为其之前值 1 和 LOWLINK[5] 之间的较小值，产生 1。

然后，由于已经考虑了所有由 1 引出的边，且 LOWLINK[1] = 1，发现 1 为强连通分支的根。该分支包含 1 和在栈中位于 1 以上的所有顶点。由于 1 为访问到的第一个顶点，顶点 2、3、4 和 5 都按序位于 1 之上。因此，栈被清空，并且将顶点 1、2、3、4、5 的列表作为 G 的一个强连通分支输出。

其他的强连通分支为 {7} 和 {6, 8}。从顶点 6 开始，计算 LOWLINK 并找出强连通分支留给读者来完成。注意，最后碰到的强连通分支的根依次为 1、7、6。 □

定理 5.4 在一个包含 n 个顶点、 e 条边的有向图 G 中，算法 5.4 能够在 $O(\max(n, e))$ 时间内正确地找到 G 的强连通分支。

证明：很容易验证，调用一次 SEARCHC(v) 所花费的时间为与顶点 v 所引出的边数成比例的时间加上一个常数，不包括对 SEARCHC 的递归调用。因此，由于在任何顶点上，仅调用 SEARCHC 一次，所以对 SEARCHC 的所有调用所需要的时间同顶点数加上边数成比例。算法 5.4 中除去 SEARCHC 的部分已经可以在 $O(n)$ 时间内实现。因此，证明了所需时间界。

为证明此证明的充分性，当 SEARCH(v) 终止时，通过对已终止 SEARCHC 调用的次数进行归纳，可正确计算 LOWLINK[v]。由 SEARCHC 的第 12~16 行， v 成为强连通分支的根，当且仅当 LOWLINK[v] = v 。进一步，按引理 5.8 的要求，输出的顶点恰好是 v 的后代，它们并不位于根在 v 之前被找到的分支中。即，在栈中 v 之上的顶点为 v 的后代，由于它们仍然处于栈中，在 v 之前没有找到它们的根。

为证明正确地计算了 LOWLINK，在图 5-15 中有两个地方需要注意，LOWLINK[v] 可以接收小于 v 的值，也就是在 SEARCHC 的第 9 和 11 行。在第一种情况下， w 为 v 的儿子，LOWLINK[w] < v 。那么，存在顶点 $x = \text{LOWLINK}[w]$ ，通过一条交叉边或回向边，从 w 的后代 y 可达。并且，含有 x 的强连通分支的根 r 为 w 的祖先。由于 $x < v$ ，有 $r < v$ ，因此 r 为 v 的某个祖先。所以，可以看到 LOWLINK[v] 至少同 LOWLINK[w] 一样小。

在第二种情形，第 11 行存在从 v 到顶点 $w < v$ 的交叉边或回向边，尚未找到其强连通分支 C 。在 C 的根 r 上所调用的 SEARCHC 还没有终止，因此 r 必定为 v 的祖先。（由于 $r \leq w < v$ ， r 要么位于 v 的左侧，要么是 v 的祖先。但是，如果 r 在 v 的左侧，SEARCHC(r) 可能不会终止。）这再次表明，LOWLINK[v] 应该至少同 w 一样小。

还必须证明，计算 LOWLINK[v] 值时 SEARCHC 是比较慢的。相反地，假定存在 v 的后代 x ，

有从 x 到 y 的交叉边或回向边, 并且含有 y 的强连通分支的根 r 为 v 的祖先。必须证明, 至少将 LOWLINK 设置为同 y 一样小。

情形 1 $x = v$ 。由归纳假设和引理 5.9, 可以假设迄今为止找到的所有强连通分支都是正确的。那么, 由于 SEARCH(r) 还未终止, y 必定还处于 STACK 中。因此, 第 11 行将 LOWLINK[v] 设置为 y 或更小的值。

情形 2 $x \neq v$ 。设 z 为 v 的儿子, x 为 v 的后代。那么, 由归纳假设, 当 SEARCHC(z) 终止时, LOWLINK[z] 已经设置为 y 或更小值。在第 9 行, LOWLINK[v] 设置为已经不能再小的值。□

5.6 路径查找问题

本节将考虑两种频繁出现的与顶点间路径有关的问题。在接下来的描述中, 设 G 为有向图。图 G^* 同 G 有相同的顶点集, 但是它含有一条从 v 到 w 的边, 当且仅当在 G 中存在一条从 v 到 w 的路径(长度为 0 或更长), 图 G^* 称为 G 的(自反和)传递闭包。

与寻找图的传递闭包紧密相关的问题是最短路径问题(shortest-path problem)。对于 G 中的每条边 e , 赋予一个非负的代价 $c(e)$ 。一条路径的代价定义为路径中各条边的代价总值。对于顶点的每个有序对 (v, w) , 最短路径问题就是找到从 v 到 w 的所有路径中的最小代价。

对于传递闭包和最短路径问题, 隐含在已知的最佳算法后面的思想是找出每对顶点之间所有路径的(无限)集合的(容易的)特殊情形。为了讨论更一般的问题, 下面引入一个特定的代数结构。

定义 闭半环(closed semiring)是一个系统 $(S, +, \cdot, 0, 1)$, 其中 S 为元素集合, $+$ 和 \cdot 为对 S 的二元操作, 满足下述 5 个性质:

1. $(S, +, 0)$ 为一个独异点(monoid), 也就是说, 在 $+$ 下是闭的[即对于 S 中所有的 a 和 b , $a + b \in S$], $+$ 是结合的[即对于 S 中的所有 a 、 b 和 c , $a + (b + c) = (a + b) + c$], 且 0 为一个单位元(identity)[即对于 S 中所有的 a , $a + 0 = 0 + a = a$]。同样地, $(S, \cdot, 1)$ 也为一个独异点。也假定 0 为一个零化子(annihilator), 即 $a \cdot 0 = 0 \cdot a = 0$ 。

2. $+$ 符合交换律(commutative), 也就是 $a + b = b + a$, 并且是幂等的(idempotent), 即 $a + a = a$ 。

3. \cdot 在 $+$ 上符合分配律(distribute), 就是说 $a \cdot (b + c) = a \cdot b + a \cdot c$, 并且 $(b + c) \cdot a = b \cdot a + c \cdot a$ 。

4. 如果 $a_1, a_2, \dots, a_i, \dots$ 为 S 中元素的可数序列, 那么 $a_1 + a_2 + \dots + a_i + \dots$ 存在并且唯一。进一步地, 结合律、交换律和幂等性可以同有限和一样适用于无限和。

5. \cdot 在可数无限和同有限和一样必定符合分配律(这一点没有遵循性质 3)。因此(4)和(5)蕴涵着

$$\left(\sum_i a_i \right) \cdot \left(\sum_j b_j \right) = \sum_{ij} a_i \cdot b_j = \sum_i \left(\sum_j (a_i \cdot b_j) \right)$$

例 5.9 下面 3 个系统为闭半环。

1. 设 $S_1 = (\{0, 1\}, +, \cdot, 0, 1)$, 其加法和乘法表如下:

$+$	0	1	\cdot	0	1
0	0	1	0	0	0
1	1	1	1	0	1

那么, 容易验证性质 1~3。对于性质 4 和 5, 注意到可数和为 0 当且仅当所有的项为 0。

2. 设 $S_2 = (R, \text{MIN}, +, +\infty, 0)$, 其中 R 为包含 $+\infty$ 的非负实数集。易验证, $+\infty$ 为 MIN 的单位元, 且 0 为 $+$ 的单位元。

3. 设 Σ 为有限字母表 (即一个符号集合), $S_3 = (F_\Sigma, \cup, \cdot, \emptyset, \{\varepsilon\})$, 其中, F_Σ 为取自 Σ 的有限长度符号串集合族, 包含空串 ε (即长度为 0 的串)。在此, 第一个运算符为集合的并, \cdot 表示集合的串联^①。 \cup 的单位元为 \emptyset , 而 \cdot 的单位元为 $\{\varepsilon\}$ 。读者可以验证性质 1~3。对于性质 4 和 5, 注意到, 如果定义 $x \in (A_1 \cup A_2 \cup \dots)$ 当且仅当对于某个 i , 有 $x \in A_i$, 则可数的集合并操作将按其方式进行。□

分析闭半环主要集中在称为闭包的一元操作, 记为 $*$ 。如果 $(S, +, \cdot, 0, 1)$ 为闭半环, 且 $a \in S$, 则定义 a^* 为 $\sum_{i=0}^{\infty} a^i$, 其中 $a^0 = 1$ 且 $a^i = a \cdot a^{i-1}$ 。也就是说, a^* 为无穷和 $1 + a + a \cdot a + a \cdot a \cdot a + \dots$ 。注意, 闭半环定义的性质 4 确保 $a^* \in S$ 。性质 4 和 5 隐含着 $a^* = 1 + a \cdot a^*$ 。注意, $0^* = 1^* = 1$ 。

例 5.10 参考例 5.9 中的半环 S_1, S_2 和 S_3 。对于 S_1 , $a^* = 1$, 其中 $a = 0$ 或 1 。对于 S_2 , 对所有在 R 中的 a , $a^* = 0$ 。对于 S_3 , 对所有 $A \in F_\Sigma$, $A^* = \{\varepsilon\} \cup \{x_1 x_2 \dots x_k \mid k \geq 1 \text{ 且 } x_i \in A, \text{ 其中 } 1 \leq i \leq k\}$ 。例如, $\{a, b\}^* = \{\varepsilon, a, b, aa, ab, ba, bb, aaa, \dots\}$, 即所有由 a 和 b 组成的串, 也包含空串。实际上, $F_\Sigma = \mathcal{P}(\Sigma^*)$, 其中, $\mathcal{P}(X)$ 表示集合 X 的幂集。□

现在, 假定有向图 $G = (V, E)$ 中, 每条边用某个闭半环 $(S, +, \cdot, 0, 1)$ 中的一个元素来标记^②。将路径标记 (label of a path) 定义为路径中边的标记的积 (\cdot), 按顺序取这些边标记。作为一种特殊情形, 零长度的路径标记为 1 (半环的 \cdot 标记)。对于每对顶点 (v, w) , 定义 $c(v, w)$ 为 v 和 w 之间所有路径标记的和。 $c(v, w)$ 将作为从 v 到 w 的代价。为了便利, 空路径集合的和为 0 (半环的 $+$ 标记)。注意, 如果 G 有环, 则在 v 和 w 之间可能存在无穷多条路径, 但是闭半环的公理保证 $c(v, w)$ 是良定义的。

例 5.11 考虑图 5-17 中的有向图, 其中每条边用例 5.9 中半环 S_1 的一个元素来标记。路径 v, w, x 的标记为 $1 \cdot 1 = 1$ 。从 w 到 w 的简单四路标记为 $1 \cdot 0 = 0$ 。实际上, 从 w 到 w 长度大于 0 的每条路径都有标记 0 。然而, 从 w 到 w 的零长度路径, 其代价为 1 。最终, $c(w, w) = 1$ 。□

现在给出一个算法为所有顶点对 v 和 w 计算 $c(v, w)$ 。算法的基本单元时间步为任意闭半环的运算符 $+$ 、 \cdot 以及 $*$ 。当然, 对于某些结构, 诸如字母表上所有字符串集组成的集合, 并不清楚能否实现这些操作, 更不用说在“单元时间”内了。然而, 对于将要使用的半环, 则很容易执行这些操作。

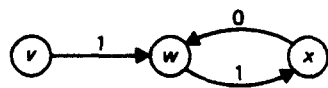


图 5-17 带标记的有向图

算法 5.5 计算顶点之间的代价。

输入: 有向图 $G = (V, E)$, 其中 $V = \{v_1, v_2, \dots, v_n\}$, 标记函数 $l: (V \times V) \rightarrow S$, 其中 $(S, +, \cdot, 0, 1)$ 为闭半环。如果 (v_i, v_j) 不在 E 中, 则取 $l(v_i, v_j) = 0$ 。

输出: 对所有在 $1 \sim n$ 之间的 i 和 j , S 中的元素 $c(v_i, v_j)$ 等于所有从 v_i 到 v_j 的路径上路径标记之和。

方法: 对所有的 $1 \leq i \leq n$, $1 \leq j \leq n$ 及 $0 \leq k \leq n$, 计算 C_{ij}^k 。其含义在于, C_{ij}^k 为从 v_i 到 v_j 的所有路径标记之和, 可能除两个端点外, 路径上的所有顶点都在集合 $\{v_1, v_2, \dots, v_k\}$ 中。例如, 路径 v_9, v_3, v_8 位于 C_{98}^4 和 C_{98}^3 中, 但是不在 C_{98}^2 中。算法如下:

① 集合 A 与 B 的串联为集合 $\{x \mid x = yz, y \in A \text{ 且 } z \in B\}$, 表示为 $A \cdot B$ 。

② 正如将在 9.1 节中要讨论的, 读者不能忽视这种情形和不确定性有穷自动机 (见 Hopcroft 和 Ullman [1969], 或者 Aho 和 Ullman [1972]) 之间的相似之处。顶点为状态, 边的标记为来自某个有穷字母表的符号。


```

begin
1.   for  $i \leftarrow 1$  until  $n$  do  $C_{ii}^0 \leftarrow 1 + l(v_i, v_i)$ ;
2.   for  $1 \leq i, j \leq n$  and  $i \neq j$  do  $C_{ij}^0 \leftarrow l(v_i, v_j)$ ;
3.   for  $k \leftarrow 1$  until  $n$  do
4.       for  $1 \leq i, j \leq n$  do
5.            $C_{ij}^k \leftarrow C_{ij}^{k-1} + C_{ik}^{k-1} \cdot (C_{kk}^{k-1})^* \cdot C_{kj}^{k-1}$ ;
6.       for  $1 \leq i, j \leq n$  do  $c(v_i, v_j) \leftarrow C_{ij}^n$ 
end

```

□

定理 5.5 算法 5.5 使用 $O(n^3)$ 个半环的 $+$ 、 \cdot 和 $*$ 操作计算 $c(v_i, v_j)$ ，其中 $1 \leq i, j \leq n$ 。

证明：易验证第 5 行执行了 n^3 次，每次需要 4 个操作，第 1、2 和 6 行的 for 循环每个重复至多 n^2 次。因此， $O(n^3)$ 个操作是足够的。

为证明正确性，对 k 归纳证明 C_{ij}^k 为所有路径标记之和，这类路径从 v_i 到 v_j 且 (除终点外) 没有索引大于 k 的中间顶点。由第 1 和 2 行，归纳的基 $k=0$ 显然，因为任何这样的路径长度都为 0，或者由唯一的边组成。由第 5 行得出归纳步，由于从 v_i 到 v_j 没有中间顶点索引大于 k 的一条路径，则

i) 没有高于 $k-1$ 的中间顶点 (项 C_{ij}^{k-1})，或者

ii) 由 v_i 指向 v_k ，再由 v_k 指向 v_j 某些次 (可能为 0 次)，最终由 v_k 指向 v_j ，整条路径上都没有高于 $k-1$ 的中间顶点 [项 $C_{ik}^{k-1} \cdot (C_{kk}^{k-1})^* \cdot C_{kj}^{k-1}$]。

闭半环定律保证算法的第 5 行正确计算这些路径的标记之和。

□

5.7 传递闭包算法

将算法 5.5 应用于两种有趣的情形。第一种为例 5.9 中描述的闭半环 S_1 。在 S_1 中，易于执行加法和乘法，且因为 $0 * = 1 * = 1$ ，执行 $*$ 也很容易。实际上，由于 1 为 \cdot 的单位元，可以使用下式取代算法 5.5 中的第 5 行^①

$$C_{ij}^k \leftarrow C_{ij}^{k-1} + C_{ik}^{k-1} \cdot C_{kj}^{k-1} \quad (5-4)$$

为了计算图的自反传递闭包，定义标记函数

$$l(v, w) = \begin{cases} 1 & \text{如果 } (v, w) \text{ 是一条边} \\ 0 & \text{如果不是} \end{cases}$$

则 $c(v, w) = 1$ 当且仅当从 v 到 w 之间存在一条长度为 0 或大于 0 的路径。使用闭半环 $\{0, 1\}$ 定律，这很容易验证。

例 5.12 考虑图 5-18 中的图，暂时忽视边上的数字。如下给定标记函数 $l(v_i, v_j)$ 。

$$\begin{array}{c|ccc}
 & v_1 & v_2 & v_3 \\
\hline
v_1 & 1 & 1 & 1 \\
v_2 & 1 & 0 & 0 \\
v_3 & 0 & 1 & 0
\end{array}$$

$k(v_i, v_j)$

然后，算法 5.5 中的第 1 行设定 $C_{11}^0 = C_{22}^0 = C_{33}^0 = 1$ 。对于 $i \neq j$ ，第 2 行设定 $C_{ij}^0 = l(v_i, v_j)$ ，因此，对于 C_{ij}^0 有下列值。

① 很自然地将这种观测解释为“仅考虑无环路径是足够的”。然而，应该注意到，用式 (5-4) 来取代第 5 行，算法 5.5 将基于路径集合来计算其和，该集合包含所有的无环路径，但也包含了一些其他的路径。

	v_1	v_2	v_3
v_1	1	1	1
v_2	1	1	0
v_3	0	1	1

C_{ij}^0

现在设定 $k=1$, 用式(5-4)取代第5行, 并执行第4和5行的循环。例如, $C_{23}^1 = C_{23}^0 + C_{21}^0 \cdot C_{13}^0 = 0 + 1 \cdot 1 = 1$ 。图 5-19 中的表给出了 C_{ij}^1 、 C_{ij}^2 和 C_{ij}^3 。

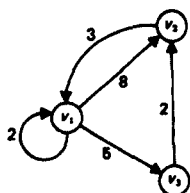


图 5-18

	v_1	v_2	v_3
v_1	1	1	1
v_2	1	1	1
v_3	0	1	1

C_{ij}^1

	v_1	v_2	v_3
v_1	1	1	1
v_2	1	1	1
v_3	1	1	1

$C_{ij}^2 = C_{ij}^1 = c(v_i, v_j)$

图 5-19 C_{ij}^k 的值

5.8 最短路径算法

为计算最短路径, 使用例 5.9 中讨论的第 2 类闭半环, 也就是包含 $+\infty$ 的非负实数。回顾一下, 实数中的加运算为 MIN, 乘运算则为加和。即要考虑结构 $(R, \text{MIN}, +, +\infty, 0)$ 。在例 5.10 中观察得到, 对于所有的 $a \in R$, $a * 0 = 0$, 因此可以删除算法 5.5 中第 5 行的 $*$ 操作, 取而代之以

$$C_{ij}^k \leftarrow \text{MIN}(C_{ij}^{k-1}, C_{ik}^{k-1} + C_{kj}^{k-1}) \quad (5-5)$$

一般地, 式(5-5)认为不经由高于 v_k 的顶点且从 v_i 到 v_j 的最短路径, 要短于

i) 不通过高于 v_{k-1} 的顶点的最短路径, 且

ii) 从 v_i 到 v_k , 再到 v_j 且尽可能短距离的路径, 在这些点之间, 没有通过高于 v_{k-1} 的顶点。

为将算法 5.5 转化为最短路径算法, 设 $l(v_i, v_j)$ 为边 (v_i, v_j) 的代价, 若边不存在则其值为 $+\infty$ 。然后用式(5-5)取代第 5 行, 并且由定理 5.5, 算法 5.5 得到的 $c(v_i, v_j)$ 值将是 v_i 和 v_j 之间所有路径中路径代价 (即边的代价之和) 最小的。

例 5.13 再次考虑图 5-18 中的图。设每条边的标记为图中表示的。图 5-20 给出了函数 l 、 C_{ij}^0 、 C_{ij}^1 、 C_{ij}^2 和 C_{ij}^3 。例如, $C_{12}^3 = \text{MIN}(C_{12}^2, C_{13}^2 + C_{32}^2) = \text{MIN}(8, 5 + 2) = 7$ 。

	v_1	v_2	v_3
v_1	2	8	5
v_2	3	∞	∞
v_3	∞	2	∞

$l(v_i, v_j)$

	v_1	v_2	v_3
v_1	0	8	5
v_2	3	0	∞
v_3	∞	2	0

C_{ij}^0

	v_1	v_2	v_3
v_1	0	8	5
v_2	3	0	8
v_3	5	2	0

C_{ij}^1

	v_1	v_2	v_3
v_1	0	8	5
v_2	3	0	8
v_3	5	2	0

C_{ij}^2

	v_1	v_2	v_3
v_1	0	7	5
v_2	3	0	8
v_3	5	2	0

$C_{ij}^3 = c(v_i, v_j)$

图 5-20 最短路径计算

5.9 路径问题与矩阵乘法

设 $(S, +, \cdot, 0, 1)$ 为闭半环。可以用通常的和与积, 定义 S 中元素的 $n \times n$ 矩阵。也就是, 设 A, B 和 C 分别包含元素 a_{ij}, b_{ij} 和 c_{ij} , 其中 $1 \leq i, j \leq n$ 。然后, 对于所有的 i 和 j , $C = A + B$ 意味着 $c_{ij} = a_{ij} + b_{ij}$, 且对所有的 i 和 j , $C = A \cdot B$ 意味着 $c_{ij} = \sum_{k=1}^n a_{ik} \cdot b_{kj}$ 。容易验证下列引理。

引理 5.10 设 $(S, +, \cdot, 0, 1)$ 为一个闭半环, 且 M_n 为 S 上 $n \times n$ 矩阵的集合。令 $+_n$ 和 \cdot_n 分别为矩阵的加和乘, 0_n 为每个元素均为 0 的矩阵, 且 I_n 为对角线上元素为 1, 其余元素为 0 的矩阵。那么, $(M_n, +_n, \cdot_n, 0_n, I_n)$ 为闭半环。

证明: 留作练习。 \square

设 $G = (V, E)$ 为一有向图, 其中 $V = (v_1, v_2, \dots, v_n)$ 。就像在算法 5.5 中一样, 假设 $l: (V \times V) \rightarrow S$ 为标记函数。设 A_c 为 $n \times n$ 矩阵, 其第 ij 元为 $l(v_i, v_j)$ 。下一个引理将通过算法 5.5 执行的计算同 S 上 $n \times n$ 矩阵的半环 M_n 中闭包 A_c^* 的计算关联起来。

引理 5.11 设 G 和 A_c 如上所述, 则 A_c^* 的第 ij 个元素为 $c(v_i, v_j)$ 。

证明: $A_c^* = \sum_{i=0}^{\infty} A_c^i$, 其中, $A_c^0 = I_n$ 且对于 $i \geq 1$, $A_c^i = A_c \cdot A_c^{i-1}$ 。通过对 i 进行归纳, 直接证得 A_c^i 的第 ij 个元素为从 v_i 到 v_j 长度为 i 的所有路径的代价之和。随后即可得到 A_c^* 的第 ij 个元素为从 v_i 到 v_j 的所有路径代价之和。 \square

作为引理 5.11 的结果, 可将算法 5.5 看作为计算矩阵闭包的算法。定理 5.5 已经阐明, 算法 5.5 需要 $O(n^3)$ 个标量运算 (scalar operation) (即取自半环的 $+$ 、 \cdot 和 $*$ 运算)。显然, 矩阵乘法的算法要求 $O(n^3)$ 个标量运算。下面证明, 当考虑取自闭半环的标量时, 两个矩阵乘积的计算在计算上等价于矩阵闭包的计算。也就是说, 给定任意算法来计算乘积, 可以构造算法来计算闭包, 反之亦然, 并且执行次数的阶也将相同。在第 6 章中更好地表达了该结果, 并证明当其标量形成一个环的时候, 少于 $O(n^3)$ 次操作就能满足矩阵的乘法。此构造分为两个部分。

定理 5.6 如果可以在 $T(n)$ 时间内计算在闭半环上任意 $n \times n$ 矩阵的闭包, 其中 $T(n)$ 满足 $T(3n) \leq 27T(n)$ [⊖], 则存在常数 c , 使两个 $n \times n$ 矩阵 A 和 B 相乘的时间 $M(n)$ 满足 $M(n) \leq cT(n)$ 。

证明: 假定要计算乘积 $A \cdot B$, A 和 B 为 $n \times n$ 矩阵。首先生成 $3n \times 3n$ 矩阵

$$X = \begin{pmatrix} 0 & A & 0 \\ 0 & 0 & B \\ 0 & 0 & 0 \end{pmatrix}$$

求出 X 的闭包。注意

$$X^2 = \begin{pmatrix} 0 & 0 & A \cdot B \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} \text{ 和 } X^3 = X^4 = \dots = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix}$$

因此

$$X^* = I_{3n} + X + X^2 = \begin{pmatrix} I_n & A & A \cdot B \\ 0 & I_n & B \\ 0 & 0 & I_n \end{pmatrix}$$

⊖ 由于最坏情况下 $T(n)$ 为 $O(n^3)$, 这是一个合理的假定。实际上, 在定理的表述中, 27 可替换为任何其他的常数。

由于可以从右上角读取 $A \cdot B$, 因此可得 $M(n) \leq T(3n) \leq 27T(n)$. \square

定理 5.6 的证明可用下面的图形解释。考虑图 G , 将其包含的 $3n$ 个顶点 $\{1, 2, \dots, 3n\}$ 划分为 3 个集合 $V_1 = \{1, 2, \dots, n\}$, $V_2 = \{n+1, n+2, \dots, 2n\}$ 和 $V_3 = \{2n+1, 2n+2, \dots, 3n\}$ 。假定 G 的每条边, 尾部顶点在 V_1 中且头部顶点在 V_2 中, 或者其尾部顶点在 V_2 中且头部顶点在 V_3 中。图 5-21 给出这样的图。设 A 和 B 为 $n \times n$ 矩阵, 其中 A 的第 ij 个元素为从顶点 i 到顶点 $n+j$ 的边的标记, 且 B 的第 ij 个元素为从顶点 $n+i$ 到顶点 $2n+j$ 的边的标记。积 $A \cdot B$ 的第 ij 个元素碰巧是从顶点 i 到顶点 $2n+j$ 的路径上标记之和, 这是由于每条这样的路径都包含从顶点 i 到某个顶点 k ($n < k \leq 2n$) 的边, 接着是一条从顶点 k 到顶点 $2n+j$ 的边。因此, 所有从 i 到 $2n+j$ 的路径上标记之和与 $A \cdot B$ 中的第 ij 个元素相同。

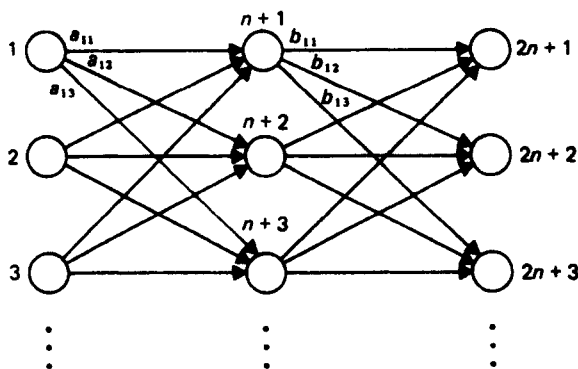


图 5-21 定理 5.6 的图形解释

现在继续证明定理 5.6 的逆命题。给定矩阵乘法算法, 忽略常数因子, 存在相同速度的闭包算法。

定理 5.7 如果能够在 $M(n)$ 时间内, 计算闭半环上任意两个 $n \times n$ 矩阵的积, 其中 $M(n)$ 满足 $M(2n) \geq 4M(n)$, 那么存在常数 c , 使得计算任意矩阵闭包的时间 $T(n)$ 满足 $T(n) \leq cM(n)$ 。

证明: 设 X 为 $n \times n$ 矩阵。首先考虑 n 为 2 的乘方, 假定为 2^k 。将 X 分解为大小为 $2^{k-1} \times 2^{k-1}$ 的 4 个矩阵,

$$X = \begin{pmatrix} A & B \\ C & D \end{pmatrix}$$

利用引理 5.11, 可以假定矩阵 X 表示图 $G = (V, E)$, 其中顶点划分为大小为 2^{k-1} 的两个集合 V_1 和 V_2 。矩阵 A 表示 V_1 中顶点之间的边, 矩阵 D 表示 V_2 中顶点之间的边。矩阵 B 表示从 V_1 中顶点到 V_2 中顶点的边, 矩阵 C 表示从 V_2 中顶点到 V_1 中顶点的边。排列如图 5-22 所示。

现在, 从 V_1 中的顶点 v_1 到同一集合中的顶点 v_2 的路径必为下述两种形式之一,

1. 永远不离开 V_1 , 或者

2. 对于某个 $k \geq 1$, $1 \leq i \leq k$, 存在顶点 w_i, x_i, y_i 和 z_i , 其中 w 和 z 在 V_1 中, x 和 y 在 V_2 中, 使得在 V_1 中有从 v_1 到 w_1 的路径, 沿着到 x_1 的边, 再沿着 V_2 中的一条路径到达 y_1 , 再沿着到达 z_1 的边, 然后沿着 V_1 中的一条路径到达 w_2 , 依此下去, 到达 z_k , 再由 V_1 中的路径到达 v_2 。图 5-23 给出其描述。

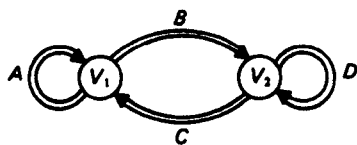


图 5-22 图 G

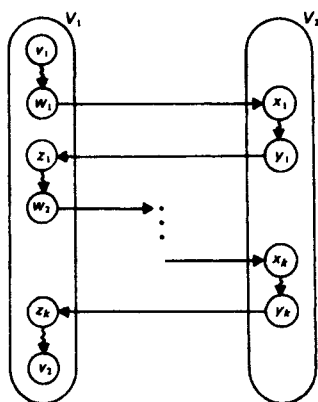


图 5-23 满足条件 2 的路径

很容易将满足条件 1 或 2 的路径的标记之和看作为 $(A + B \cdot D^* \cdot C)^*$ 。也就是说, 在图 5-23 中, $B \cdot D^* \cdot C$ 表示从 w_i 到 x_i 到 y_i 再到 z_i 的路径。因此, $A + B \cdot D^* \cdot C$ 表示 V_1 中顶点之间的路径, 这些路径只是一条边, 或者直接转到 V_2 , 停在 V_2 中, 再回到 V_1 中。可以将 V_1 中顶点之间的所有路径连续表示为由 $A + B \cdot D^* \cdot C$ 表示的路径。因此, $(A + B \cdot D^* \cdot C)^*$ 表示 V_1 中顶点之间的所有路径。因而, X^* 的左上部分包含 $(A + B \cdot D^* \cdot C)^*$ 。

设 $E = (A + B \cdot D^* \cdot C)^*$ 。出于相似的原因, 可以将 X^* 的 4 个部分写为:

$$X^* = \begin{pmatrix} E & E \cdot B \cdot D^* \\ D^* \cdot C \cdot E & D^* + D^* \cdot C \cdot E \cdot B \cdot D^* \end{pmatrix} = \begin{pmatrix} E & F \\ G & H \end{pmatrix}$$

可以通过下列步骤, 计算这 4 个部分 E 、 F 、 G 和 H :

$$T_1 = D^*$$

$$T_2 = B \cdot T_1$$

$$E = (A + T_2 \cdot C)^*$$

$$F = E \cdot T_2$$

$$T_3 = T_1 \cdot C$$

$$G = T_3 \cdot E$$

$$H = T_1 + G \cdot T_2$$

上述步骤需要 $2^{k-1} \times 2^{k-1}$ 矩阵的 2 个闭包、6 次乘法以及 2 次加法。因此, 有:

$$T(1) = 1 \quad (5-6)$$

$$T(2^k) \leq 2T(2^{k-1}) + 6M(2^{k-1}) + 2 \cdot 2^{2k-2}, \quad k > 1$$

式(5-6)右边的 3 项分别表示闭包、乘法和加法的代价。要求存在常数 c , 对所有的 k , $T(2^k) \leq cM(2^k)$ 满足式(5-6)。由于可以将 c 选为任意大, 归纳基础 $k=0$ 时是显然的。在归纳步, 假定对某个 c 有 $T(2^{k-1}) \leq cM(2^{k-1})$ 。从定理的假设 $M(2n) \geq 4M(n)$, 有 $M(2^{k-1}) \geq 2^{2k-2}$ 。[注意, $M(1)=1$]。因此, 从式(5-6)

$$T(2^k) \leq (2c+8)M(2^{k-1})$$

再次, 由定理的假设, $M(2^{k-1}) \leq 1/4M(2^k)$, 有

$$T(2^k) \leq (c/2+2)M(2^k) \quad (5-7)$$

如果选择 $c \geq 4$, 式(5-7)意味着 $T(2^k) \leq cM(2^k)$, 正是要证明的。

如果 n 不是 2 的乘方, 可以将 X 嵌入到维数为 2 的乘方的大矩阵中, 其形式为

$$\begin{pmatrix} X & 0 \\ 0 & I \end{pmatrix}$$

其中 I 为具有最小可能大小的单位方阵。这将使矩阵的大小至多为原来的两倍, 并因此通过 8 的因子来增加常量。因此, 存在新的常数 c' , 无论 n 是否为 2 的乘方, 对所有的 n , 有 $T(n) \leq c'M(n)$ 。

□

推论 1 计算布尔矩阵的闭包所需要的时间等同于计算相同大小的两个布尔矩阵的乘积所需要的时间。

第 6 章将看到计算布尔矩阵乘积的渐近快速算法, 因此表明, 可以在少于 $O(n^3)$ 的时间内计算图的传递闭包。

推论 2 用操作 MIN 和 + 起到标量加和乘的作用, 计算非负实数矩阵的闭包所需要的时间与计算两个该类型矩阵的乘积所耗费的时间同阶。

到目前为止, 对于所有顶点对的最短路径问题, 没有已知的少于 cn^3 时间的方法, 其中 c 是常数。

5.10 单源问题

在多数应用中, 可能仅要找到从一个顶点(源)开始的最短路径。实际上, 可能期望仅找到两个特定顶点之间的最短路径, 但是较之最佳单源算法, 对此问题在最坏情况下还没有已知的比其更有效的算法。

对于单源问题, 尚未有同闭半环和算法 5.5 相似的概念。并且, 如果仅了解存在一条从源开始的路径的顶点, 则问题很简单, 能够通过许多算法, 在 e 条边的图中, 在 $O(e)$ 时间内完成。由于没有“考察”过所有边的算法不可能是正确的, 不难相信, $O(e)$ 是所能期望的最好时间了。

当提及找到单源最短路径时, 情况有所变化。虽然假定算法应该需要超过 $O(e)$ 时间是什么理由的, 但是不超过 $O(e)$ 时间的算法并无人所知。将讨论一个 $O(n^2)$ 的算法, 构造一个顶点集 S , 这些顶点从源开始的最短距离是已知的。在每一步, 将剩余顶点 v 加入 S , 该顶点从源开始的距离是剩余顶点中最短的。如果所有边都有非负的代价, 则确保从源到 v 的路径仅通过 S 中的顶点。因此, 对于每个顶点 v , 仅需要记录从源开始, 经由一条仅通过 S 中顶点的路径到 v 的最短距离。现在形式化地给出其算法。

算法 5.6 单源最短路径算法(Dijkstra 算法)。

输入: 有向图 $G=(V, E)$ 、源 $v_0 \in V$ 以及从边映射到非负实数的函数 l 。如果 (v_i, v_j) 不是一条边, $v_i \neq v_j$, 则取 $l(v_i, v_j)$ 为 $+\infty$, 并且 $l(v, v)=0$ 。

输出: 对每个 $v \in V$, 所有从 v_0 到 v 的路径 P 中边的标记之和的最小值。

方法: 构造集合 $S \subseteq V$, 使得从源到 S 中每个顶点 v 的最短路径整个位于 S 中。数组 $D[v]$ 包含从 v_0 到 v 仅经由 S 中顶点的当前最短路径的代价。图 5-24 给出该算法。

□

```

begin
1.  S ← {v0};
2.  D[v0] ← 0;
3.  for each v in V - {v0} do D[v] ← l(v0, v);
4.  while S ≠ V do
      begin
5.      choose a vertex w in V - S such that D[w] is a minimum;
6.      add w to S;
7.      for each v in V - S do
8.      D[v] ← MIN(D[v], D[w] + l(w, v))
      end
end
end

```

图 5-24 Dijkstra 算法

例 5.14 考虑图 5-25。首先, $S = \{v_0\}$, $D[v_0] = 0$, 对于 $i = 1, 2, 3, 4$, $D[v_i]$ 分别为 $2, +\infty, +\infty, 10$ 。在第 4~8 行循环中的第一次迭代时, 由于 $D[v_1] = 2$ 为 D 中的最小值, 选择 $w = v_1$ 。然后, 设定 $D[v_2] = \min(+\infty, 2+3) = 5$, $D[v_4] = \min(10, 2+7) = 9$ 。图 5-26 中概括性地给出了 D 中的值序列和算法 5.6 中的其他计算。

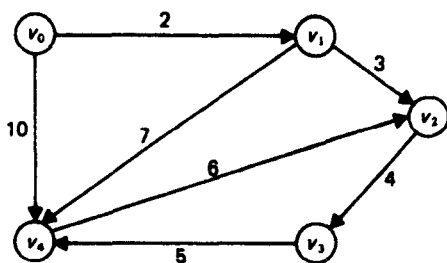


图 5-25 带标记边的图

Iteration	S	w	$D[w]$	$D[v_1]$	$D[v_2]$	$D[v_3]$	$D[v_4]$
Initial	$\{v_0\}$	—	—	2	$+\infty$	$+\infty$	10
1	$\{v_0, v_1\}$	v_1	2	2	5	$+\infty$	9
2	$\{v_0, v_1, v_2\}$	v_2	5	2	5	9	9
3	$\{v_0, v_1, v_2, v_3\}$	v_3	9	2	5	9	9
4	All	v_4	9	2	5	9	9

图 5-26 基于图 5-25, 算法 5.6 的计算

定理 5.8 算法 5.6 计算从 v_0 到每个顶点的最短路径的代价, 需要 $O(n^2)$ 时间。

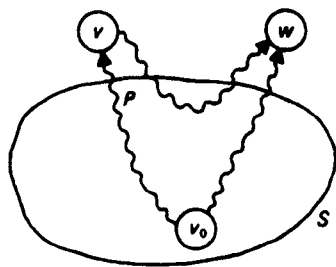
证明: 同在第 5 行中选择 w 一样, 第 7~8 行中的 for 循环要求 $O(n)$ 步。这些是第 4~8 行中 while 循环的主要代价。而后者需要执行 $O(n)$ 次, 其总代价为 $O(n^2)$ 。显然, 第 1~3 行需要 $O(n)$ 时间。

对于正确性, 必须基于 S 的大小进行归纳证明, 对于 S 中的每个 v , $D[v]$ 等于从 v_0 到 v 的一条最短路径的长度。而且, 对所有的 $v \in V - S$, $D[v]$ 为从 v_0 到 v 且除 v 本身外, 路径整体位于 S 中的最短路径的长度。

归纳基础。 $\|S\| = 1$ 。从 v_0 到其本身的最短路径的长度为 0, 从 v_0 到 v 且除 v 本身外, 路径整体位于 S 中的路径含有唯一边 (v_0, v) 。因此, 第 2 和 3 行正确地初始化数组 D 。

归纳步。假定在第 5 行选择顶点 w 。如果 $D[w]$ 不是从 v_0 到 w 最短路径的长度, 那么必然存在一条最短路径 P 。除 w 外, 路径 P 必定包含不在 S 中的顶点。设 v 为 P 上第一个这样的顶点。但是另一方面, 从 v_0 到 v 的距离却比 $D[w]$ 短, 并且, 到 v 的最短路径除 v 本身外, 整体位于 S 中。(参见图 5-27)。因此, 由归纳假设, 当选择 w 时, 有 $D[v] < D[w]$, 矛盾。得出结论, 不存在路径 P , 且 $D[w]$ 为从 v_0 到 w 的最短路径的长度。

由第 8 行, 归纳假设的第二部分, 即 $D[w]$ 保持正确性, 显而易见。

图 5-27 到顶点 v 的路径

5.11 有向无环图的支配集：概念整合

本章已经说明了设计有效图算法的几种技术，诸如深度优先搜索和合理计算顺序的安排等。第4章研究了许多有用的数据结构，用于表示由各种操作来操纵集合。本章将通过一个例子说明，结合第4章中的数据结构和本章的图技术可以设计出非常有效的算法。具体就是使用离线 MIN 算法、不相交集合并算法和 2-3 树，结合本章的深度优先搜索技术，为有向无环图寻找支配集 (dominators)。

如果有向图 $G = (V, E)$ 的顶点 r 到 V 中每个顶点都存在一条路径，则以顶点 r 为根。在本节的余下部分，将假定 $G = (V, E)$ 为以 r 为根的带根有向无环图。

如果从根到 w 的每条路径都包含 v ，则顶点 v 为 w 的支配点 (dominator)。每个顶点是其本身的支配点，并且根支配着每个顶点。顶点 w 的支配点的集合，可依其在从根到 w 的最短路径上出现的顺序线性排序。将最靠近 w 的 w 的支配点 (除 w 本身之外) 称为 w 的直接支配点 (immediate dominator)。由于每个顶点的支配点线性排列，可用一棵含有根 r 的树表示关系“ v 支配 w ”，该树称为 G 的支配树 (dominator tree)。支配集的计算能用于代码优化问题 (见 [Aho 和 Ullman, 1973] 和 [Hecht, 1973])。

对一棵有 e 条边的带根有向无环图，为计算其支配树，需要开发一个 $O(e \log e)$ 算法。该算法主要是为了说明如何将这里的技术同前一章结合起来。该算法基于下述 3 个引理。

设 $G = (V, E)$ 为一个图， $G' = (V', E')$ 为 G 的子图。如果 $(a, b) \in E'$ ，写为 $a \Rightarrow_{G'} b$ 。分别用 $\vec{a} \Rightarrow_{G'} b$ 和 $\overset{\leftarrow}{a} \Rightarrow_{G'} b$ 表示 $\Rightarrow_{G'}$ 的传递闭包和自反传递闭包。如果 $(a, b) \in E$ ，写作 $a \Rightarrow b$ 。

引理 5.12 设 $G = (V, E)$ 为带根 r 的有向无环图， $S = (V, T)$ 为带根 r 的图 G 的深度优先生成树。设 a, b, c 和 d 为 V 中的顶点，满足 $a \Rightarrow b \Rightarrow c \Rightarrow d$ 。设 (a, c) 和 (b, d) 为前向边。然后，用前向边 (a, d) 取代前向边 (b, d) 不会改变 G 中任何顶点的支配点 (见图 5-28)。

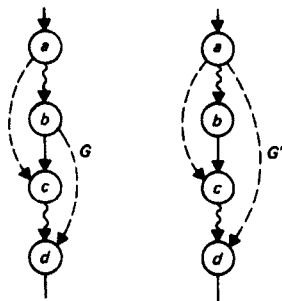


图 5-28 引理 5.12 的转化

证明：用边 (a, d) 取代 G 中的边 (b, d) ，设 G' 为由此得到的图。假定 v 为 w 的支配点，它在 G 中而不在 G' 中。那么，在 G' 中存在一条从 r 到 w 的路径 P ，不包含 v 。显然，由于边 (a, d) 是只在 G' 中而不在 G 中的边，路径 P 必须使用这条边。因此，可以写 $P: r \Rightarrow a \Rightarrow d \Rightarrow w$ 。这使得 v 必须位于路径 $a \Rightarrow d$ 上，并且 v 相异于 a 和 d 。如

果在深度优先计数中， $a < v \leq b$ ，那么用路径 $a \Rightarrow c \Rightarrow d$ 来取代 P 中的边 $a \Rightarrow d$ ，在 G 中得到一条从 r 到 w ，不含 v 的路径。如果 $b < v < d$ ，那么用路径 $a \Rightarrow b \Rightarrow d$ 来取代 P 中的边 $a \Rightarrow d$ ，在 G 中得到一条从 r 到 w 且不含 v 的路径。由于 $a < v \leq b$ 或者 $b < v < d$ ，所以在 G 中存在一条从 r 到 w 且不含 v 的路径。矛盾。

假定 v 为 w 的关节点，在 G' 中但不在 G 中。那么，在 G 中存在一条从 r 到 w 且不含 v 的路径 P 。由于该路径不在 G' 中，路径 P 必定包含边 $b \Rightarrow d$ 。这使 v 位于路径 $b \Rightarrow d$ 上，并且 $b < v < d$ 。因此，在 G' 中存在一条路径 $r \Rightarrow a \Rightarrow d \Rightarrow v$ 不含 v ，矛盾。 □

引理 5.13 设 G 和 S 与引理 5.12 中的定义相同。设 $a \Rightarrow b$ ，并且 (a, b) 为 G 中的前向边。如果不存在使 $c < a$ 且 $a < d \leq b$ 的前向边 (c, d) ，也不存在使 $a < d \leq b$ 的交叉边 (e, d) ，那么，删除指向 b 的树边不会改变一个顶点的支配集 (见图 5-29)。

证明：证明相似于引理 5.12。 □

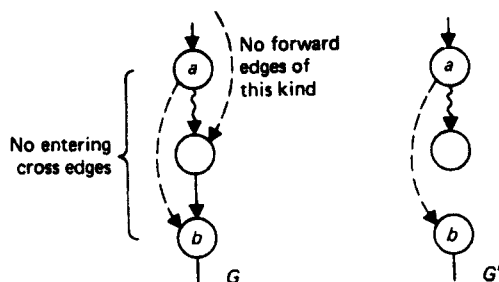


图 5-29 引理 5.13 的转化

引理 5.14 设 G 和 S 与引理 5.12 中相同, (c, d) 为 G 的一条交叉边, b 为 c 和 d 的最高深度数的共同祖先。设 a 为 $\text{MIN}(\{b\} \cup \{v \mid (v, w) \text{ 为一条前向边且 } b < w \leq c\})$ 。如果, 除顶点 b 外, 不存在进入路径 $b \xrightarrow{S} c$ 上任何顶点的交叉边, 那么用前向边 (a, d) 取代交叉边 (c, d) 不会改变 G 中任何顶点的支配集 (见图 5-30)。

证明: 用边 (a, d) 取代 G 中的边 (c, d) , G' 为由此得到的图。假定顶点 v 为顶点 w 的支配点, 在 G 中但不在 G' 中。那么, G 中存在一条从 r 到 w 且不含 v 的路径 P 。显然, P 必定包含边 (a, d) 。因此, v 必定在生成树路径 $a \xrightarrow{S} b$ 上, 否则用 $a \xrightarrow{S} d$ 或者 $a \xrightarrow{S} c \Rightarrow d$ 取代 $a \Rightarrow d$, 由此在 G 中得到从 r 到 w 不含 v 的路径。然后, 用 $a \Rightarrow u \xrightarrow{S} c \Rightarrow d$ 取代 P 中的边 $a \Rightarrow d$, 其中 u 在路径 $b \xrightarrow{S} c$ 上且 $u > b$, 由此在 G 中得到一条从 r 到 w 不含 v 的路径, 矛盾。

假定 v 为 w 的支配点, 在 G' 中但不在 G 中。那么在 G 中存在一条从 r 到 w 不含 v 的路径 P 。显然, P 包含边 (c, d) 。将 P 写为 $r \xrightarrow{S} c \Rightarrow d \xrightarrow{S} w$ 。由于没有头端位于 $b \xrightarrow{S} c$ 上 (b 除外) 的交叉边, 路径 $r \xrightarrow{S} c$ 必定包含路径 $a \xrightarrow{S} b$ 上的某个顶点。设 x 为最高深度数的这类顶点。设 P_1 为路径 P 中从 r 到 x 的一部分, 其后有 $x \xrightarrow{S} d$, 接着是 P 中从 d 到 w 的部分。设 P_2 为路径 $r \xrightarrow{S} a \Rightarrow d$, 其后有 P 中从 d 到 w 的部分。如果 v 位于 P_1 , 那么 v 必定在 $x \xrightarrow{S} d$ 上, 且 $v > x$ 。如果 v 位于 P_2 , 那么 v 必定在 $r \xrightarrow{S} a$ 上。由于 $a \leq x$, G' 中这些路径中的一条不包含 v , 矛盾。 \square

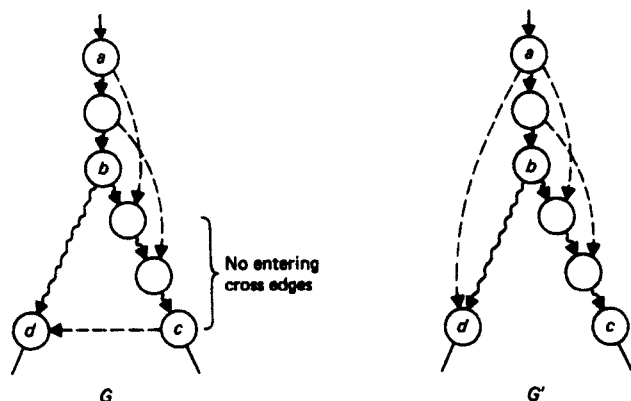


图 5-30 引理 5.14 中的转化

显然, 重复应用引理 5.12 ~ 5.14 中的转化, 直到它们不再可用, 将把 G 转化为一棵树。由

于转化不改变支配关系, 最终的树必定是 G 的支配树。这将是计算 G 的支配集的算法。整个技术就是设计一个数据结构, 允许应用引理 5.12、5.13 和 5.14 的转化有效地找到合适的边。

直观地, 按如下方式为给定的有向无环图 $G = (V, E)$ 构造支配树。首先, 从根开始执行 G 的深度优先搜索, 构造一棵深度优先生成树 $S = (V, T)$ 。根据其深度优先数, 重新标记 G 的顶点。然后, 将引理 5.12 ~ 5.14 中的转化应用于 G 。用两个相互关联的程序实现转化, 一个用于处理前向边, 另一个用于交叉边。

1. 前向边

假定 G 不存在交叉边。如果顶点 v 作为多于一条前向边的头端, 那么, 由引理 5.12, 除了边尾端最小的那条外, 所有以 v 为头的前向边都可以删除, 而不会改变任意顶点的支配集。

组合边 (composite edge) 是一个有序对 (t, H) , 其中 t 为称为组合边尾的顶点, H 为组合边头的顶点集合。组合边 $(t, \{h_1, h_2, \dots, h_k\})$ 表示边 $(t, h_1), (t, h_2), \dots, (t, h_k)$ 的集合。

重复应用引理 5.12 改变各前向边的尾端。有时, 对于有相同尾端 t 的边集合中的每条边, 其尾端将变为 t' 。为了更有效, 用一条组合边表示具有共同尾端的前向边的某些集合。初始时, 每条前向边 (t, h) 用组合边 $(t, \{h\})$ 表示。

将 G 的每个顶点 v 关联到集合 $F[v]$ 。 $F[v]$ 包含形如 $(t, \{h_1, h_2, \dots, h_m\})$ 的对, 使 t 为 v 的祖先, 每个 h_i 为 v 的后代, 并且 (t, h_i) 为 G 中的一条前向边。首先, $F[v] = \{(t, \{v\})\}$, 其中 t 为以 v 为头的一条前向边的最小深度数的尾端。

现在, 按逆先序遍历生成树。当沿着生成树边 (v, w) 返回时, 发现对于 w 的每个确定的祖先 t , 集合 $F[w]$ 包含一条组合边, t 为当前 w 的后代的一条前向边的尾端。然后, 执行下述操作。

1. 设 $(t, \{h_1, h_2, \dots, h_m\})$ 为 $F[w]$ 中的组合边, 包含最大深度数的尾端。如果 $t = v$, 则从 $F[w]$ 中删除该组合边 (组合边表示前向边的集合, 应用引理 5.12, 这些前向边的尾端已经位于 v 之前, 但不会更靠前)。对于 $1 \leq i \leq m$, 从 G 中删除以 h_i 为头的生成树边 (这一步对应引理 5.13 的应用)。

2. 若 G 中存在前向边 (t, v) , \ominus 则对于使得 $u \geq t$ 的 $F[w]$ 中的每条组合边 $(u, \{h_1, \dots, h_m\})$, 做如下操作:

a) 从 $F[w]$ 中删除 $(u, \{h_1, \dots, h_m\})$ 。

b) 将 $\{h_1, \dots, h_m\}$ 与代表其他边的 $F[v]$ 中组合边的头与边 (t, v) 相合并。

(该步相应于引理 5.12 的应用。)

3. 用 $F[v] \cup F[w]$ 取代 $F[v]$ 。

例 5.15 考虑图 5-31 所示的带根有向无环图 G 。对 G 进行深度优先搜索, 生成图 5-32a 所示的图。在图中也给出了与每个顶点相关的 F 集合。图 5-32b ~ d 给出了前向边处理的结果。图 5-32d 中是最后得到的支配树。 □

当到达根的时候得到的图是正确的支配树 (假定没有交叉边), 这留给读者来证明。将不相交集合并算法用来处理组合边边头的集合, 可以有效地实现该算法。由于必然可以有效地删除一条组合边, 可以用 2-3 树表示由组合边组成的集合 $F[v]$, 使得在组合边集合中找到拥有最大深度数尾端的组合边, 得到组合边集合的并。用这样的数据结构处理含有 e 条边的图需要的时间为 $O(e \log e)$ 。

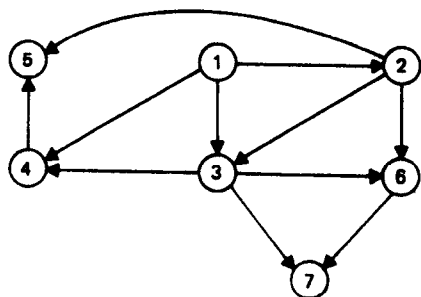


图 5-31 带根有向无环图

\ominus 假定除了其尾端深度数最小的前向边之外, 已经从 G 中删除了所有以 v 为头的前向边。

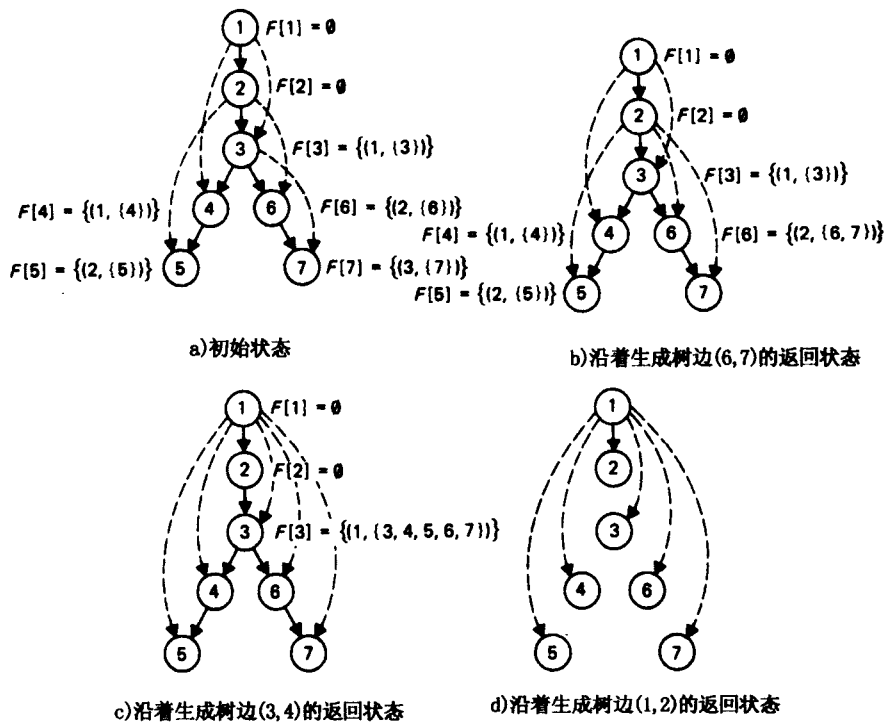


图 5-32 前向边转化的结果

II. 交叉边

一般地,不能假定不存在交叉边。通过以下将要描述的方法,可以用前向边取代交叉边。然而,由于在 I 中为了有效地应用引理 5.14 而构建了数据结构,因而,在如 I 那样处理前向边之前,不应该做这种取代。并且,由于去除的每条交叉边都将变为一条前向边,所以在去除交叉边之前不应该完全地执行 I 。应该在关于前向边倒转的先序遍历过程中添加交叉边的处理步。注意,由于使用了引理 5.13, I 要求在某些时刻,不存在指向某些顶点的交叉边。实际上,用下面所描绘的步骤,按反先序进行的遍历应该使读者相信,在交叉边的存在使得不能使用引理 5.13 之前,每条交叉边都将已经改变为一条前向边。

设 S 为 G 的深度优先生成树。首先,对于每条交叉边 (v, w) , 计算 v 和 w 深度数最大的共同祖先。对每个顶点 v , 使其对应由有序对 (u, w) 组成的集合 $L[v]$, 其中 $u > w$, (u, w) 表示对 u 和 w 深度数较高祖先的一次请求。初始时, $L[v] = \{(v, w) \mid \text{存在交叉边}(v, w), v > w\}$ 。在如同 I 部分那样遍历 S 时, 执行下述操作。

1. 遍历树边 (v, w) , $v < w$ 时, 从 $L[v]$ 中删除满足 $y \geq w$ 的每个 (x, y) 。顶点 v 为 x 和 y 的具有最大深度数的共同祖先。
2. 由生成树边 (v, w) 返回到 v 的时候, 用 $L[v] \cup L[w]$ 取代 $L[v]$ 。

通过一般化习题 4.21 中离线 MIN 算法, 可以在至多 $O(eG(e))$ 步内计算深度数最大的祖先, 其中, e 为 G 中的边数。

通过应用引理 5.14 处理交叉边, 将其转化为前向边。当处理前向边时, 必须将交叉边转化到前向边。对于每个顶点 v , 将组合边的集合 $C[v]$ 关联。起初, $C[v] = \{(v, \{h_1, \dots, h_m\}) \mid \text{对于 } 1 \leq i \leq m, (v, h_i) \text{ 为一条交叉边}\}$ 。当通过树边 (v, w) 返回到顶点 v 的时候, 除了处理前向

边的步骤之外, 执行下述步骤。

1. 用 $C[v] \cup C[w]$ 取代 $C[v]$;
2. 删除交叉边 (x, y) , 则在当前表示 v 的组合边中, v 为 x 和 y 最大深度数的祖先。如果组合边有尾端 t , 用前向边 (t, y) 取代交叉边 (x, y) 。如果已经存在一条指向 y 的前向边, 则仅保留有最小尾端的前向边;
3. 如果存在的话, 设 (u, v) 为以 v 为头的前向边。否则, 设 (u, v) 为指向 v 的树边。在查找了 v 的所有后代之后, 从 $C[v]$ 删除尾端位于 u 之上的组合交叉边。将已删除的组合交叉边的头端的集合合并在一起, 生成一条以 u 为尾端的新组合边, 将其添加到 $C[v]$ 。

例 5.16 考虑图 5-33a 所示的深度优先生成树。对于选定的顶点给出了 $C[v]$ 值。由于存在一条从顶点 2 指向顶点 8 的前向边, 将 $C[8]$ 中的组合边 $(8, \{4\})$ 变为 $(2, \{4\})$ 。然后将 $C[8]$ 合并到 $C[7]$ 。由于 $(1, 7)$ 为一条前向边, 将组合边 $(2, \{4\})$ 变为 $(1, \{4\})$ 。当回到 6 时, $C[6]$ 变成为 $\{(1, \{4\}), (6, \{5\})\}$ 。

当从顶点 6 回到顶点 3 时, $C[3]$ 变为 $\{(3, \{5\}), (1, \{4\})\}$ 。顶点 3 为顶点 6 和 5 以及顶点 8 和 4 的最大深度数的祖先, 因此, 从 $C[3]$ 中删除组合边 $(3, \{5\})$ 和 $(1, \{4\})$, 并将前向边 $(3, 5)$ 和 $(1, 4)$ 添加到 G 中。结果如图 5-33b 所示。其余查找没有引起进一步的变化。 \square

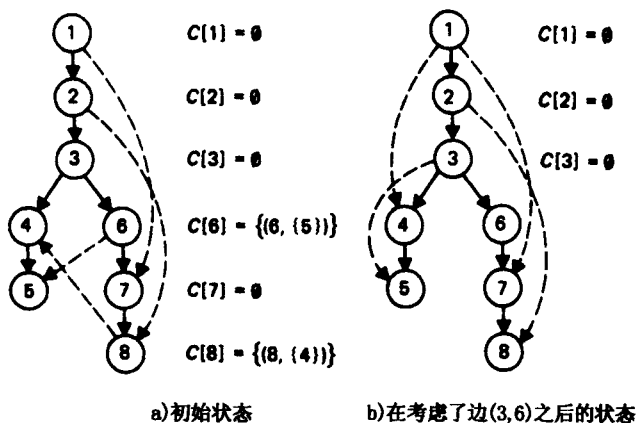


图 5-33 删除交叉边

可以用 2-3 树表示组合交叉边, 支配算法的正式描述作为一个有趣的习题留给读者。如果能够将会合适的数据结构组合起来, 那么就掌握了第 4 章和第 5 章的技术。

习题

- 5.1 对于图 5-34 所示, 假定所给出的边是无向的, 找出其最小代价生成树。
- 5.2 设 $S = (V, T)$ 为算法 5.1 构造的一棵最小代价生成树。 T 中各条边的代价为 $c_1 \leq c_2 \leq \dots \leq c_n$ 。设 S' 为任意一棵生成树, 其各边的代价为 $d_1 \leq d_2 \leq \dots \leq d_n$ 。证明: 对于 $1 \leq i \leq n$, $c_i \leq d_i$ 。
- **5.3 使用下述策略, 在 $O(e)$ 时间内为有 n 个顶点和 e 条边的图找出一棵最小代价生成树, e 和 n 满足 $e \geq f(n)$, 其中函数 f 须由你给出。在各个时刻, 将顶点分组为集合, 并将这些集合用目前所找到的树边连接。把所有以集合中 1 或 2 个顶点为入射点的边保留在集合的优先队列中。初始时, 每个顶点自身位于一个集合中。

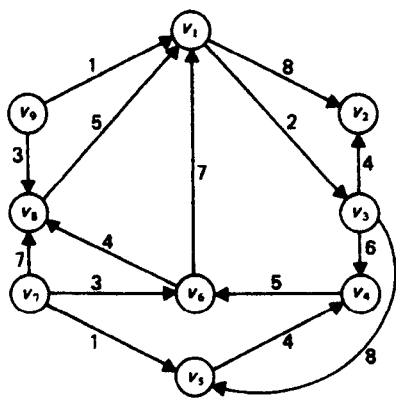


图 5-34

迭代步是为了找出最小集合 S , 以及导出 S 的最小代价边, 假定为集合 T 。然后, 将那条边加到生成树中, 合并集合 S 和 T 。然而, 如果所有的集合大小至少为 $g(n)$, 其中 g 为必须找出的另一个函数, 那么对每个集合用一个顶点构造一个新图。如果初始时, S_1 中的某个顶点同 S_2 中的某个顶点是邻接的, 那么在新图中, 与集合 S_1 和 S_2 相关的顶点是邻接的。新图中, 连接 S_1 和 S_2 边的代价是初始时 S_1 中任意顶点和 S_2 中某个顶点之间任意边中的最小代价。然后, 将算法递归地应用于新图。

需要解决的问题是选择 $g(n)$, 使得 $f(n)$ 最小。

- 5.4 为图 5-35 中的无向图找出深度优先生成森林。对于每棵树, 随意选择初始顶点。找出图的双连通分支。

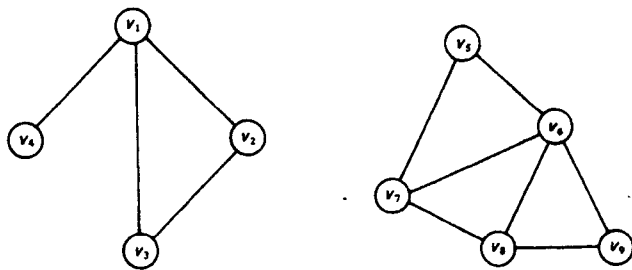


图 5-35

- 5.5 为图 5-34 中的有向图找出深度优先生成森林。然后找出强连通分支。

- 5.6 为图 5-36 找出双连通分支。

- 5.7 借助深度优先搜索设计有效算法, 实现下列任务:

- 将无向图分解为其连通分支;
- 在一个无向连通图中找到一条路径, 使其在每个方向恰好通过每条边一次;
- 验证是否一个无向图无环;
- 为一个无环有向图中的顶点找到一个顺序, 如果存在一条从 v 到 w 且长度大于 0 的路径, 则 $v < w$;
- 确定是否一个连通无向图的边可以有向化, 得到一个强连通有向图。[提示: 证明目标当且仅当从 G 中删除

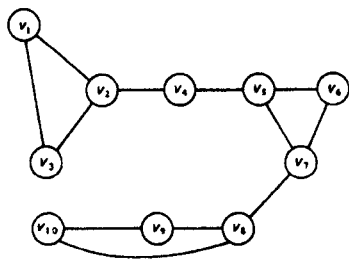


图 5-36

任意边得到的是一个连通图时才能实现。]

- 5.8 设 $G=(V, E)$ 为多重图 (multigraph) (即无向图 G 中, 每对顶点之间不止有一条边)。编写一个 $O(\|E\|)$ 算法删除度为 2 的顶点 v [通过用边 (u, w) 取代边 (u, v) 和 (v, w)], 并用一条边取代边的多个副本。注意: 通过一条边取代边的多个副本可能会创建度为 2 的顶点, 而必须删除这类顶点。相似地, 删除度为 2 的顶点可能创建多重边, 随后必须删除这类边。
- 5.9 无向图的欧拉回路是一条路径, 起始并终止于相同的顶点, 并且经过每条边恰好一次。一个连通无向图 G 包含一个欧拉回路, 当且仅当每个顶点的度数都为偶数。给出一个 $O(e)$ 算法, 在包含 e 条边的图中找出一个欧拉回路, 如果存在的话。
- *5.10 设 $G=(V, E)$ 为双连通无向图, (v, w) 为 G 中的一条边, $G'=(\{v, w\}, \{(v, w)\})$ 。找出一种技术在线执行形如 $\text{FINDPATH}(s, t)$ 的指令序列, 其中 s 为 G' 的一个顶点, (s, t) 为在 G 中而不在 G' 中的一条边。 $\text{FINDPATH}(s, t)$ 的执行包括在 G 中找到一条起始于边 (s, t) , 终止于 G' 中的顶点而不是 s 的简单路径, 并且将路径中的顶点和边加入 G' 中。任何序列的执行时间都应该为 $O(\|E\|)$ 。
- 5.11 考虑图 5-37 中的有向图 G 。
- 找出 G 的传递闭包。
 - 找出 G 中每对顶点间最短路径的长度。图 5-37 给出每条边的代价。
- *5.12 找出一个含有 4 个元素的闭半环。
- *5.13 有向图 $G=(V, E)$ 的传递归约 (transitive reduction) 定义为边数尽可能少的图 $G'=(V, E')$, 使得 G' 的传递闭包等价于 G 的传递闭包。证明: 如果 G 无环, 则 G 的传递归约唯一。
- **5.14 证明: 当 (合理地) 假定 $8R(n) \geq R(2n) \geq 4R(n)$ 且 $8T(n) \geq T(2n) \geq 4T(n)$ 时, 为一个 n 顶点的无环图计算传递归约的时间 $R(n)$ 为计算传递闭包的时间 $T(n)$ 的常量因子。
- *5.15 证明: 当假设同习题 5.14 时, 为一个无环图计算传递归约所需要的时间, 对任意图找出其传递归约所需要时间的常量因子。
- *5.16 对于传递闭包, 证明习题 5.15。
- *5.17 设 A 为闭半环 $\{0, 1\}$ 上的 $n \times n$ 矩阵。不使用图形化解释, 证明下述结论:
- $A^* = I_n + A + A^2 + \cdots + A^n$
 - $\begin{pmatrix} A & B \\ 0 & C \end{pmatrix}^* = \begin{pmatrix} A^* & A^* B C^* \\ 0 & C^* \end{pmatrix}$
- [提示: 证明
- $$\begin{pmatrix} A & B \\ 0 & C \end{pmatrix}^i = \begin{pmatrix} A^i & A^{i-1} B + A^{i-2} B C + \cdots + B C^{i-1} \\ 0 & C^i \end{pmatrix}]$$
- *5.18 证明: 如果某些边的代价为负, 但不包含代价为负的环, 则 5.8 节中的最短路径算法仍然适用。如果存在代价为负的环, 则结果会怎样?
- **5.19 证明: 含有 $+\infty$ 和 $-\infty$ 的正负实数不是一个闭半环。在此条件下, 如何解释习题 5.18?

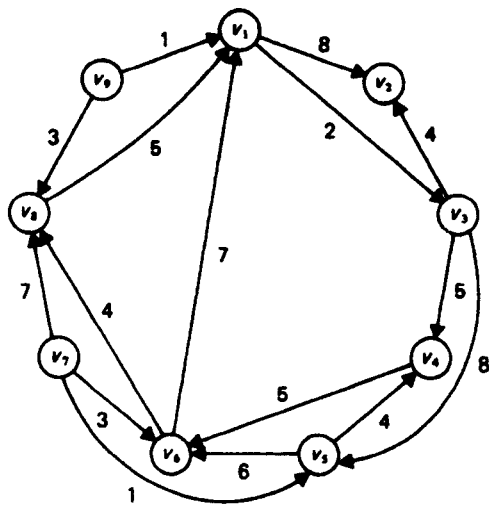


图 5-37

[提示: 在算法 5.5 中实际使用到的是闭半环的什么属性?]

- 5.20 使用算法 5.6, 找出从顶点 v_0 到图 5-37 的图 G 中每个顶点 v 的最短距离。
- *5.21 证明: 对于包含 e 条边和 n 个顶点的图, 可以在 $O(e \log n)$ 时间内解决非负边的单源最短路径问题。[提示: 使用合适的数据结构, 当 $e < n^2$ 的时候, 可以有效地执行算法 5.6 的第 5 和 8 行。]
- **5.22 证明: 在包含 e 条边和 n 个顶点的图中, 对于任意固定的常数 k , 可以在 $O(ke + kn^{1+1/k})$ 时间内解决非负边的单源最短路径问题。
- *5.23 证明: 对于含有负代价边而不包含负代价环的任意图, 可以用图 5-38 中的单源最短路径算法计算从 v_0 到每个顶点 v 的最短路径。

```

begin
  S ← {v0};
  D[v0] ← 0;
  for each v in V - {v0} do D[v] ← l(v0, v);
  while S ≠ V do
    begin
      choose a vertex w in V - S such that D[w] is a minimum;
      S ← S ∪ {w};
      for v ∈ V such that D[v] > D[w] + k(w, v) do
        begin
          D[v] = D[w] + l(w, v);
          S ← S - {v}
        end
      end
    end
  end
end

```

图 5-38 单源最短路径算法

- *5.24 对于图 5-38 中的算法, 其执行时间的阶是多少? [提示: 在一个包含 5 个顶点的图中执行该算法, 其中边的代价如图 5-39 所示。]

	1	2	3	4	5
1	0	7	8	9	10
2	0	0	8	9	10
3	0	-2	0	9	10
4	0	-4	-3	0	10
5	0	-7	-6	-5	0

图 5-39 一个 5 顶点图边的代价

- 5.25 给出一个算法, 确定对于含有正和负代价边的有向图是否含有代价为负环。
- 5.26 在图 5-38 的算法中, 改变 w 的选择规则, 对任意代价边, 保证时间界为 $O(n^3)$ 。
- 5.27 给定一个正整数的 $n \times n$ 矩阵 M , 编写一个算法找出起始于 $M[n, 1]$, 终止于 $M[1, n]$ 的邻接矩阵元序列, 使得邻接矩阵元之间差的绝对值之和最小。如果 (a) $i = k \pm 1$ 且 $j = l$, (b) $i = k$ 且 $j = l \pm 1$, 则两个矩阵元 $M[i, j]$ 和 $M[k, l]$ 是邻接的。例如, 在图 5-40 中, 序列 7, 5, 8, 7, 9, 6, 12 为一种解决方案。
- 5.28 对于每个顶点 $v \in V$, 从 v_0 到 v 的所有路径 P , 图 5-24 中的算法计算图 5-40 正整数矩阵路径 P 代价的最小值。对于 V 中的每个顶点 v , 修改算法以得到一条最小代价的路径。
- **5.29 编写一个程序, 实现 5.11 节中的支配集算法。

1	9	6	12
8	7	3	5
5	9	11	4
7	3	2	6

研究性问题

- 5.30 存在众多图问题, 必然有一些应用深度优先搜索技术。其中之一所表示的是关于 k 连通性的问题。对于无向图中的每对顶点 v 和 w , 如果在 v 和 w 之间存在 k 条路径, 使得(除了 v 和 w 外)没有顶点出现在一条以上路径上, 则该无向图是 k 连通的。因此, 双连通指的是2连通的。Hopcroft和Tarjan[1973b]给出了一个线性时间算法找出3连通分支。这很自然地联想到, 对于每个 k , 存在线性(在大量顶点和边中)时间算法找到 k 连通分支。你能够找到一个吗?
- 5.31 对于最小代价生成树的一个有趣期望是设计一个线性(与边数相关)时间的算法, 可以涉及也可以不涉及深度优先搜索技术。
- 5.32 值得考虑的第3个问题是当 $e < n^2$ 时的单源最短路径问题。存在 $O(e)$ 算法来找到两个特定点之间的最短距离吗? 读者应该从Johnson[1973]中了解到习题5.21和5.22, 也应该从Spira和Pan[1973]中了解到对于那些仅使用边代价和之间比较的算法, 证明其一般需要 $n^2/4$ 次比较。在边的权重为小整数的情况下, Wagner[1974]使用桶排序技术得到了一个 $O(e)$ 算法。
- 5.33 已经证明, 如果仅允许操作MIN和+, 则找出所有点对之间最短路径的问题需要 kn^3 步, 其中常数 $k > 0$ (Kerr, [1972])。M. O. Rabin将这一结果精化为 $n^3/6$ 。然而, 如果允许其他操作, 可能我们能够做得比 $O(n^3)$ 更好。例如, 如果使用布尔操作之外的其他操作, 可以在少于 $O(n^3)$ 步内完成传递闭包(等价地, 布尔矩阵乘法), 下一章将看到这类结果。

文献与注释

与图论相关的两种资源为Berge[1958]和Harary[1969]。关于最小代价生成树的算法5.1来自Kruskal[1956]。Prim[1957]给出了该问题的另一种方法。P. M. Spira提供了习题5.3中提及的算法。

使用深度优先搜索的双连通性算法由J. E. Hopcroft提出。而强连通分支算法则来自Tarjan[1972]。在文献中, 许多应用通过使用深度优先搜索为已知算法得到优化。Hopcroft和Tarjan[1973a]给出一个线性的平面测试算法。Hopcroft和Tarjan[1973b]描述了一个3连通分支的线性算法。Tarjan[1973a]利用这种概念得到了目前为止找出支配集的最优算法, 并且Tarjan[1973b]给出了一个“流图可归约性”的测试算法。

作为通用的路径查找算法, 算法5.5基本上归功于Kleene[1956], 他将其同“正则表达式”(参见9.1节)联系起来。这里所给出的算法形式来源于McNaughton和Yamada[1960]。 $O(n^3)$ 的传递闭包算法归功于Warshall[1962]。相似地, 所有点对的最短路径算法来自Floyd[1962]。(也见Hu[1968])。Dijkstra[1959]提出了单源算法。Spira和Pan[1973]证明了在决策树模型下, Dijkstra算法必定是优的。

在Dantzig、Blattner和Rao[1967]中研究发现, 如果不存在负环, 则负边的存在不会影响所有点的最短路径问题。Johnson[1973]讨论了含有负边的单源问题, 在其中可以找到习题5.21和5.22的解答。为了找到最短路径, Spira[1973]给出了一个 $O(n^2 \log^2 n)$ 的期望时间算法。

路径问题和矩阵乘法问题之间的关系源自Munro[1971]和Furman[1970](定理5.7), 以及Fischer和Meyer[1971](定理5.6)。与传递归约有关的关系(习题5.13~5.15)来自Aho、Garey和Ullman[1972]。Even[1973]讨论了图的 k 连通性问题。

第6章 矩阵乘法及相关操作

本章研究元素取自任意环的矩阵乘法的渐近计算复杂度。我们将看到“普通” $O(n^3)$ 矩阵乘法算法是可以渐近改进的， $O(n^{2.81})$ 的复杂度足以解决两个 $n \times n$ 矩阵的乘法运算。同时，本章还将研究可以转化为矩阵乘法的其他运算，比如 LUP 分解、矩阵求逆、行列式计算等，它们与矩阵乘法同阶。可以看到，矩阵乘法可以转化为矩阵求逆运算，因此，这两种运算中的一种的渐近时间提高都将改善另一种运算。最后，本章以两个低于 $O(n^3)$ 渐近时间复杂度的布尔矩阵乘法算法作为总结。

对于本章所描述的大规模算法，目前的计算机硬件实际情况难以满足。首先，由于潜在的常数因子，需要足够大的 n 值才能使算法真正渐近优于普通的 $O(n^3)$ 。其次，对这类算法的数值误差的控制目前还没有深入的理解。但是，通过本章内容的学习，在处理此类算法问题时，读者可以认识到平时显而易见的算法并不一定是最优的，寻找更有效或更切实际的算法来解决这类问题总是需要的。

6.1 基础知识

本节主要给出矩阵乘法所涉及到的代数概念的基本定义。如果读者对环及线性代数知识比较熟悉，可以直接进入 6.2 节。

定义 环是代数结构 $(S, +, \cdot, 0, 1)$ ，其中 S 是元素的集合， $+$ 和 \cdot 分别是集合 S 的二元操作。 S 中的任意元素 a, b 和 c 满足以下性质：^①

1. $(a+b)+c=a+(b+c)$ 和 $(a \cdot b) \cdot c = a \cdot (b \cdot c)$ ($+$ 和 \cdot 符合结合律)
2. $a+b=b+a$ ($+$ 符合交换律)
3. $(a+b) \cdot c = a \cdot c + b \cdot c$ 和 $a \cdot (b+c) = a \cdot b + a \cdot c$ (\cdot 对 $+$ 符合分配律)
4. $a+0=0+a=a$ (0 是 $+$ 的单位元)
5. $a \cdot 1=1 \cdot a=a$ (1 是 \cdot 的单位元)
6. 对于 S 中的任意元素 a ，存在逆 $-a$ ，使得 $a+(-a)=(-a)+a=0$

注意，在上述最后一个属性中，加法逆的存在不一定适用于任意的闭半环（见 5.6 节）。同样，闭半环的第 4 条性质即无限和存在且值唯一，也不总是适合于环。如果 \cdot 符合交换律，则该环称为交换环。

如果在交换环中，任意元素 a 存在乘法逆 a^{-1} ，使得 $a \cdot a^{-1} = a^{-1} \cdot a = 1$ ，则该环为域。

例 6.1 基于普通加法 $+$ 和乘法 \cdot 的实数集形成环，但不形成闭半环。

系统 $(\{0, 1\}, +, \cdot, 0, 1)$ ，其中 $+$ 是模 2 的加法运算， \cdot 是普通的乘法运算，形成环，但不是闭半环，因为 $1+1+\dots$ 没有按照半环的性质进行良定义。如果 $+$ 重新定义如下：如果 a 和 b 是 0，则 $a+b$ 是 0；否则 $a+b$ 是 1。在这种情况下，产生例 5.1 的闭半环 S_2 。 S_2 不是环，因为 1 没有逆。□

下面的定义及引理中将介绍一类重要的由矩阵生成的环。

定义 假设 $R=(S, +, \cdot, 0, 1)$ 为环， M_n 为元素来自 R 的 $n \times n$ 矩阵的集合， 0_n 为 $n \times n$

① 这里采用 $M[i, j]$ 表示矩阵 M 第 i 行第 j 列的元素。

的零矩阵, I_n 为 $n \times n$ 的单位矩阵, 其主对角线元素为 1, 其他为 0。对于 M_n 中任意矩阵 A 和 B , 设 $A +_n B$ 为 $n \times n$ 的矩阵 C , 其中 $C[i, j] = A[i, j] + B[i, j]$, $A \cdot_n B$ 为 $n \times n$ 的矩阵 D , 其中 $D[i, j] = \sum_{k=1}^n A[i, k] \cdot B[k, j]$ 。

引理 6.1 $(M_n, +_n, \cdot_n, 0_n, I_n)$ 是环。

证明: 基本练习。 \square

注意, 尽管在基础环 R 中乘法运算符合交换律, 但上面定义的矩阵乘法运算 \cdot_n 在 $n > 1$ 时不符合交换律。然而, 在不引起混淆的情况下, 我们将用基础环中的加法和乘法运算符号 $+$ 和 \cdot 来代替 $+_n$ 和 \cdot_n , 另外, 经常略去明显存在的乘法运算符号。

假定 R 为环, M_n 为 $n \times n$ 矩阵环, 其中矩阵里的元素来自 R 。假定 n 是偶数, 则在 M_n 中的 $n \times n$ 矩阵可以分成 4 个 $(n/2) \times (n/2)$ 矩阵。假定 $R_{2,n/2}$ 为 2×2 矩阵环, 矩阵里的元素来自 $M_{n/2}$ 。可以明显看出, 在 M_n 中的 $n \times n$ 矩阵乘法及加法运算与 $R_{2,n/2}$ 中的 2×2 矩阵的乘法及加法运算分别是等价的, 其中 2×2 矩阵的元素为 $(n/2) \times (n/2)$ 矩阵。

引理 6.2 假定 f 为 M_n 到 $R_{2,n/2}$ 的映射, 使得 $f(A)$ 是 $R_{2,n/2}$ 中的如下矩阵:

$$\begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}$$

其中 A_1 为 A 的左上 $1/4$ 象限, A_2 为 A 的右上 $1/4$ 象限, A_3 为 A 的左下 $1/4$ 象限, A_4 为 A 的右下 $1/4$ 象限。有以下性质:

$$\text{i)} f(A+B) = f(A) + f(B)$$

$$\text{ii)} f(AB) = f(A) \cdot f(B)$$

证明: 这是用 $R_{2,n/2}$ 中的 $+$ 和 \cdot 定义代替 $M_{n/2}$ 中的 $+$ 和 \cdot 定义的基本练习。 \square

引理 6.2 的重要性在于能够通过 2×2 及 $(n/2) \times (n/2)$ 的矩阵乘法算法来构造 $n \times n$ 矩阵的乘法算法。在下一节中将通过这种方法构造渐近更快速的矩阵乘法算法。在这里需要注意的是这里的 2×2 矩阵乘法不是任意的 2 阶方阵的乘法, 而是针对 $R_{2,n/2}$ 专门设计的矩阵乘法。因为 $R_{2,n/2}$ 中的乘法并不要求满足交换律, 即使 R 满足, 2×2 的矩阵乘法算法也不能假定 $(n/2) \times (n/2)$ 乘法操作符合交换律。当然, 可以使用环的任何性质。

定义 假定 A 为矩阵元素来自某个域的 $n \times n$ 矩阵, 如果存在 A 的逆, 记为 A^{-1} , 使得 $AA^{-1} = I_n$ 成立, 则 A^{-1} 也是 $n \times n$ 矩阵。

可以看出, 如果 A^{-1} 存在, 则必唯一, 且 $AA^{-1} = A^{-1}A = I_n$ 成立; 而且有 $(AB)^{-1} = B^{-1}A^{-1}$ 成立。

定义 假定 A 为任意的 $n \times n$ 矩阵, A 的行列式表示为 $\det(A)$, 是下列积的整数 $1 \sim n$ 的所有排列 $p = (i_1, i_2, \dots, i_n)$ 的和。

$$(-1)^k \prod_{j=1}^n A[j, i_j]$$

其中, 如果 p 是偶排列, [从 $(1, 2, \dots, n)$ 能够偶数次交换] 则 k_p 为 0; 如果 p 是奇排列, [从 $(1, 2, \dots, n)$ 通过奇数次交换得到] 则 k_p 为 1。

很容易证明, 一个排列是偶排列, 或者是奇排列, 但不会两者都是。

例 6.2 令

$$A = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix}$$

1, 2, 3 的 6 个排列为 $(1, 2, 3)$, $(1, 3, 2)$, $(2, 1, 3)$, $(2, 3, 1)$, $(3, 1, 2)$ 和 $(3, 2, 1)$ 。 $(1, 2, 3)$ 是偶排列, 因为它只需要交换 0 次就能实现。排列 $(2, 3, 1)$ 也是偶排列, 因

为可以先从(1, 2, 3)开始交换 2 和 3 的位置得到(1, 3, 2), 然后再交换 1 和 2 得到(2, 3, 1)。类似地, (3, 1, 2)是偶排列, 其他 3 个排列为奇排列, 可得

$$\det(A) = a_{11}a_{22}a_{33} - a_{11}a_{23}a_{32} - a_{12}a_{21}a_{33} + a_{12}a_{23}a_{31} + a_{13}a_{21}a_{32} - a_{13}a_{22}a_{31} \quad \square$$

假定 A 为元素来自某个域的 $n \times n$ 矩阵, 证明 A^{-1} 存在当且仅当 $\det(A) \neq 0$, 同时, $\det(AB) = \det(A)\det(B)$ 。如果 $\det(A) \neq 0$, 称矩阵 A 为非奇异阵。

定义 $m \times n$ 矩阵 A 称为上三角矩阵, 如果对于任意的 $1 \leq j < i \leq m$, 有 $A[i, j] = 0$ 。 $m \times n$ 矩阵 A 称为下三角矩阵, 如果对于任意的 $1 \leq i < j \leq n$, 有 $A[i, j] = 0$ 。

两个上三角矩阵的积还是上三角矩阵, 同样两个下三角矩阵的乘法运算结果也是下三角矩阵。

引理 6.3 如果 A 是上三角或下三角矩阵, 则

a) $\det(A)$ 为矩阵主对角元素的乘积 (即 $\prod_i A[i, i]$);

b) A 是非奇异矩阵当且仅当主对角线上没有 0 元素。

证明: (a) 除 $(1, 2, \dots, n)$ 之外的任意排列 (i_1, i_2, \dots, i_n) 至少存在两个数 i_k 和 i_j , 使得 $i_j < j$ 且 $i_k > k$, 因此, $\det(A)$ 中的所有项都是 0, 除了排列 $(1, 2, \dots, n)$ 所对应的项; (b) 由 (a) 直接得证。 \square

定义 单位矩阵 (unit matrix) 是主对角线元素全是 1 的矩阵 (非主对角线元素值任意)。

由观察可知, 单位上三角矩阵和单位下三角矩阵的行列式是单位矩阵。

定义 如果在矩阵中, 任意一行或者一列都有且仅有一个 1, 其他元素为 0, 则称为置换阵。

定义 在矩阵 A 中删除某些行和某些列后所剩下的矩阵称为 A 的子阵。 $n \times n$ 矩阵 A 的主子阵为矩阵 A 的前 k 行和前 k 列构成的方子阵, 其中 $1 \leq k \leq n$ 。

矩阵 A 的秩记为 $\text{rank}(A)$, 表示矩阵 A 的最大非奇异方子阵的大小。例如 $n \times n$ 置换阵的秩为 n 。

观察 $A = BC$ 的情况, 满足 $\text{rank}(A) \leq \min(\text{rank}(B), \text{rank}(C))$ 。如果 A 有 m 行, 且它的秩为 m , 则任意从 A 中取 k 行所形成的矩阵的秩为 k 。

定义 矩阵 A 的转置记为 A^T , 是通过把所有的 i 和 j 交换 $A[i, j]$ 和 $A[j, i]$ 的位置形成的矩阵。

6.2 Strassen 矩阵乘法算法

假设 A 和 B 为两个 $n \times n$ 矩阵, 其中 n 为 2 的幂。通过引理 6.2, 对矩阵 A 和 B 进行划分, 每个矩阵里面划分成 4 个 $(n/2) \times (n/2)$ 的矩阵, 用这些 $(n/2) \times (n/2)$ 的矩阵表示 A 和 B 的积为以下形式:

$$\begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \begin{bmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{bmatrix} = \begin{bmatrix} C_{11} & C_{12} \\ C_{21} & C_{22} \end{bmatrix}$$

其中

$$\begin{aligned} C_{11} &= A_{11}B_{11} + A_{12}B_{21} \\ C_{12} &= A_{11}B_{12} + A_{12}B_{22} \\ C_{21} &= A_{21}B_{11} + A_{22}B_{21} \\ C_{22} &= A_{21}B_{12} + A_{22}B_{22} \end{aligned} \quad (6-1)$$

如果把 A 和 B 看成 2×2 矩阵, 矩阵的每个元素为 $(n/2) \times (n/2)$ 矩阵, 那么矩阵 A 和 B 的乘法结果可以按照式 (6-1) 表示的 $(n/2) \times (n/2)$ 个矩阵的和与积进行计算。假定计算 C_{ij} 需要 $(n/2) \times (n/2)$ 矩阵的 m 次乘法操作和 a 次加 (或减法) 操作。采用递归算法,

可以在 $T(n)$ 时间内计算两个 $n \times n$ 矩阵的乘法, 其中 n 是 2 的幂, 则 $T(n)$ 满足以下不等式:

$$T(n) \leq mT\left(\frac{n}{2}\right) + \frac{an^2}{4}, \quad n > 2 \quad (6-2)$$

在式(6-2)右边的第一项为进行 m 对 $(n/2) \times (n/2)$ 矩阵乘法所需的代价, 第二项为进行 a 次加法所需要的代价, 假定 $n^2/4$ 为每对 $(n/2) \times (n/2)$ 矩阵进行加或减运算所需要的时间。通过与定理 2.1 类似的分析方法, 在 $c=2$ 的情况下, 可以证明只要 $m > 4$, 存在常数 k , 式(6-2)的解的上限为 $kn^{\log m}$ 。解的形式不受 a 的影响。在 $m < 8$ 的情况下, 可以得到渐近复杂度优于普通 $O(n^3)$ 的算法。

V. Strassen 发现了一个只采用 7 次乘法运算就能计算出两个 2×2 矩阵的乘积的巧妙算法, 矩阵的元素来自任意的环。通过递归使用该方法, 可以在 $O(n^{\log 7})$ 时间内完成两个 $n \times n$ 的矩阵乘法, 阶大约在 $n^{2.81}$ 左右。

引理 6.4 两个 2×2 矩阵的乘法可以通过 7 次乘法运算和 18 次加/减运算得到结果, 其中矩阵的元素选自任意环。

证明: 计算以下矩阵的乘法

$$\begin{bmatrix} c_{11} & c_{12} \\ c_{21} & c_{22} \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \begin{bmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{bmatrix}$$

首先计算以下积。

$$m_1 = (a_{12} - a_{22})(b_{21} + b_{22})$$

$$m_2 = (a_{11} + a_{22})(b_{11} + b_{22})$$

$$m_3 = (a_{11} - a_{21})(b_{11} + b_{12})$$

$$m_4 = (a_{11} + a_{12})b_{22}$$

$$m_5 = a_{11}(b_{12} - b_{22})$$

$$m_6 = a_{22}(b_{21} - b_{11})$$

$$m_7 = (a_{21} + a_{22})b_{11}$$

然后采用以下公式计算 c_{ij} :

$$c_{11} = m_1 + m_2 - m_4 + m_6$$

$$c_{12} = m_4 + m_5$$

$$c_{21} = m_6 + m_7$$

$$c_{22} = m_2 - m_3 + m_5 - m_7$$

计算次数是显然的。而所期望的 c_{ij} 值的证明是基于环的定律的简单代数练习。□

习题 6.5 给出了一种只需要 7 次乘法和 15 次加法运算就可计算两个 2×2 矩阵乘法的方法。

定理 6.1 矩阵元素来自任意环的两个 $n \times n$ 矩阵乘法运算可以在 $O(n^{\log 7})$ 内完成。

证明: 首先, 考虑 $n=2^k$ 的情况, 令 $T(n)$ 为两个 $n \times n$ 矩阵相乘所需的算术操作次数。运用引理 6.4 可得,

$$T(n) = 7T\left(\frac{n}{2}\right) + 18\left(\frac{n}{2}\right)^2 \quad \text{对于 } n \geq 2$$

因此, 通过简单修改定理 2.1, 可得 $T(n)$ 为 $O(7^{\log n})$ 或 $O(n^{\log 7})$ 。

在 n 不是 2 的幂的情况下, 可以把该矩阵嵌入到维数最小的比 2 的幂更高的矩阵中。在这种情况下, 矩阵的维至多增加一倍, 因此, 常数因子至多增加 7。因此, 对所有 $n \geq 1$ 的情况,

$T(n)$ 为 $O(n^{\log 7})$ 成立。□

定理 6.1 只考虑 $T(n)$ 的函数增长率。但要确定 n 值为何时 Strassen 算法的性能优于普通的算法, 则必须确定比例常数的值。因此, 在 n 不是 2 的幂的情况下, 仅把矩阵嵌入到维数最小的比 2 的指数幂更高的矩阵中将会产生一个大的常数。在这里把矩阵嵌入到维为 $2^p r$ 的矩阵中, 其中 r 为适当小的值, 运用 p 次引理 6.4, 采用 $O(r^3)$ 的算法计算两个 $r \times r$ 矩阵的乘法。因此, 可以提出更为通用的递归算法, 当 n 为偶数时, 像前面一样, 把每个矩阵划分成四个子矩阵, 当 n 为奇数时, 先把维数增加 1, 然后采用上述算法。

需要指出, 确定比例常数仅涉及到算术运算的次数。要比较 Strassen 算法与普通的矩阵乘法算法, 还需考虑访问矩阵元素所需的额外复杂度。

6.3 矩阵求逆

在这节中, 将说明一类矩阵求逆运算可以简化为矩阵乘法的问题。这类矩阵包括所有可以求逆的三角矩阵, 但不包含所有的可逆矩阵。在后面, 将把该方法推广到任意的可逆矩阵。在本节和 6.4、6.5 节中, 假定所有矩阵的元素选自某个域。

引理 6.5 假设矩阵 A 划分如下:

$$\begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}$$

假设 A_{11}^{-1} 存在。定义 $\Delta = A_{22} - A_{21}A_{11}^{-1}A_{12}$, 且假定 Δ^{-1} 存在, 则

$$A^{-1} = \begin{bmatrix} A_{11}^{-1} + A_{11}^{-1}A_{12}\Delta^{-1}A_{21}A_{11}^{-1} & -A_{11}^{-1}A_{12}\Delta^{-1} \\ -\Delta^{-1}A_{21}A_{11}^{-1} & \Delta^{-1} \end{bmatrix} \quad (6-3)$$

证明: 通过直接的算术操作可以直接得到

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} = \begin{bmatrix} I & 0 \\ A_{21}A_{11}^{-1} & I \end{bmatrix} \begin{bmatrix} A_{11} & 0 \\ 0 & \Delta \end{bmatrix} \begin{bmatrix} I & A_{11}^{-1}A_{12} \\ 0 & I \end{bmatrix}$$

其中 $\Delta = A_{22} - A_{21}A_{11}^{-1}A_{12}$, 可得

$$\begin{aligned} A^{-1} &= \begin{bmatrix} I & -A_{11}^{-1}A_{12} \\ 0 & I \end{bmatrix} \begin{bmatrix} A_{11}^{-1} & 0 \\ 0 & \Delta^{-1} \end{bmatrix} \begin{bmatrix} I & 0 \\ -A_{21}A_{11}^{-1} & I \end{bmatrix} \\ &= \begin{bmatrix} A_{11}^{-1} + A_{11}^{-1}A_{12}\Delta^{-1}A_{21}A_{11}^{-1} & -A_{11}^{-1}A_{12}\Delta^{-1} \\ -\Delta^{-1}A_{21}A_{11}^{-1} & \Delta^{-1} \end{bmatrix} \end{aligned} \quad \square$$

引理 6.5 并不适合所有的非奇异矩阵。例如, $n \times n$ 置换矩阵 A 的元素 $A[i, j]$ 在 $j = n - i + 1$ 时, 等于 1; 否则等于 0, 矩阵 A 是个非奇异矩阵, 但对于任意主子阵 A_{11} , 行列式 $\det(A_{11}) = 0$ 。引理 6.5 适合所有的非奇异上三角或者下三角矩阵。

引理 6.6 如果 A 为非奇异的上(下)三角矩阵, 则引理 6.5 中的矩阵 A_{11} 和 Δ 可逆, 且逆矩阵为非奇异的上(下)三角矩阵。

证明: 假定 A 为上三角矩阵。下三角矩阵的证明方法类似。显然, A_{11} 为非奇异的上三角矩阵, 因此 A_{11}^{-1} 存在。观察可得 $A_{21} = 0$, 因此 $\Delta = A_{22} - A_{21}A_{11}^{-1}A_{12} = A_{22}$, 可得 Δ 为非奇异的上三角矩阵。□

定理 6.2 令 $M(n)$ 为环上的两个 $n \times n$ 矩阵做乘法运算所需的时间, 假设对所有的 m , $8M(m) \geq M(2m) \geq 4M(m)$ 成立, 则存在一个常数 c , 使得任意的 $n \times n$ 的非奇异上(下)三角矩阵 A 的逆可以在 $cM(n)$ 时间内计算。

证明: 这里只证明 n 是 2 的幂的情况。很显然, 在 n 不是 2 的幂的情况下可以把矩阵 A 嵌入到如下形式的矩阵中:

$$\begin{bmatrix} A & 0 \\ 0 & I_m \end{bmatrix}$$

其中, $m+n \leq 2n$ 为 2 的幂。通过增加常数 c , 使之至多为 8 的因子, 这个结论可以适合任意的 n 。^①

假定 n 是 2 的幂, 可以把矩阵 A 划分成 4 个 $(n/2) \times (n/2)$ 子矩阵, 递归地使用式 (6-3)。 $A_{21} = 0$, 所以 $\Delta = A_{22}$ 。因此, 需要 $2T(n/2)$ 时间求三角矩阵 A_{11} 和 Δ 的逆, $2M(n/2)$ 时间用于两个非平凡阵的乘法, $n^2/4$ 时间用于右上三角矩阵的求反运算。从定理假定及观察 $M(1) \geq 1$, 有 $n^2/4 \leq M(n/2)$ 。因此

$$\begin{aligned} T(1) &= 1 \\ T(n) &\leq 2T\left(\frac{n}{2}\right) + 3M\left(\frac{n}{2}\right) \quad \text{对于 } n \geq 2 \end{aligned} \quad (6-4)$$

证明式 (6-4) 蕴涵着 $T(n) \leq 3M(n)/2$ 是个比较基础的练习。 \square

6.4 矩阵的 LUP 分解

LUP 分解是求解线性方程组的有效方法。

定义 对于 $m \leq n$ 的 $m \times n$ 矩阵 A , 可分解成一对矩阵 L 和 U , 使得 $A = LU$, 其中 L 是 $m \times m$ 的单位下三角矩阵, U 是 $m \times n$ 的上三角矩阵, L 和 U 称为矩阵 A 的 LU 分解。

可以解方程 $Ax = b$, 求 x , 其中 A 为 $n \times n$ 矩阵, x 为 n 维的未知列向量, b 为 n 维列向量, 把矩阵 A 分解成单位下三角矩阵 L 和上三角矩阵 U 的乘积, 假定这两个因式存在。则方程 $Ax = b$ 可以写为 $LUx = b$ 。为了解 x , 可以先解 $Ly = b$, 求出 y , 然后由 $Ux = y$ 解出 x 。

该方法难点在于矩阵 A 可能不存在 LU 分解, 即使 A 为非奇异矩阵, 如矩阵:

$$\begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

为非奇异矩阵, 却没有 LU 分解。如果 A 为非奇异矩阵, 则存在置换矩阵 P , 使得 AP^{-1} 存在 LU 分解。在此给出一种算法, 对于任意非奇异矩阵 A , 寻找因式 L , U 和 P , 使得 $A = LUP$ 。矩阵 L , U 和 P 称为矩阵 A 的 LUP 分解。

算法 6.1 LUP 分解。

输入: $n \times n$ 非奇异矩阵 M , 其中 n 为 2 的幂。

输出: 矩阵 L , U 和 P , 使得 $M = LUP$, 其中 L 为单位下三角矩阵, U 为上三角矩阵, P 为置换矩阵。

方法: 调用 FACTOR(M, n, n), FACTOR 为图 6-4 所示的递归过程。FACTOR 的算法过程如图 6-1、6-2 和 6-3 所示, 在图中, 矩阵的阴影部分表示其中的元素全部为 0。

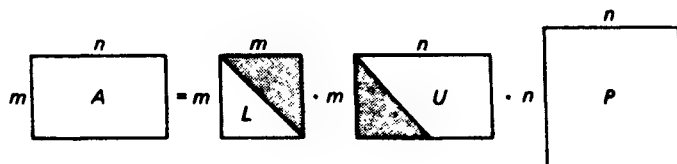


图 6-1 FACTOR 的期望结果

① 进一步的分析可得到适合任意整数 n 的常数 c , 但与 n 只是 2 的幂的情况所得到的最好结果有较大的区别。

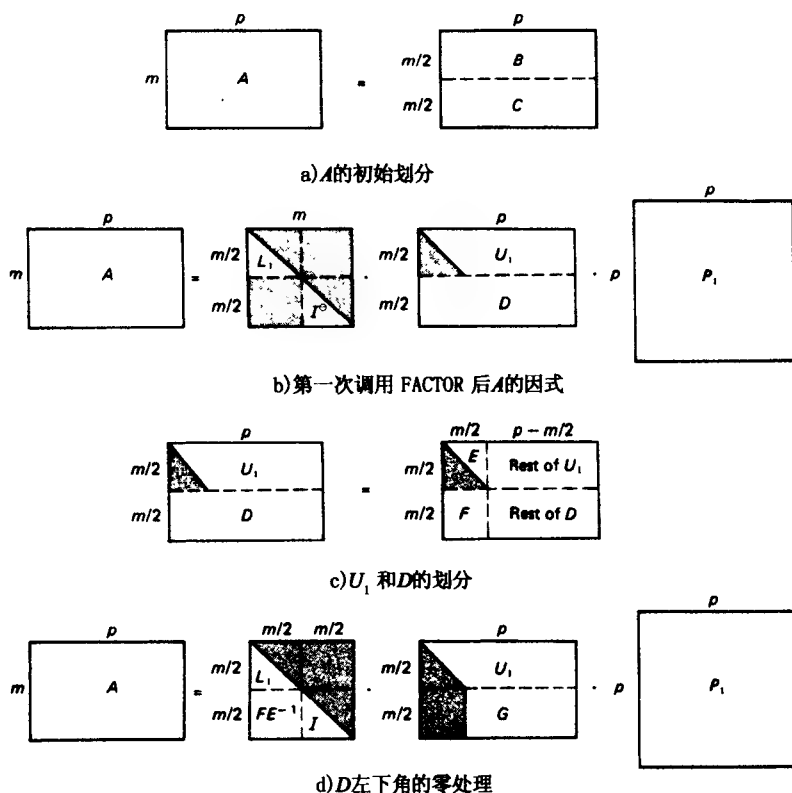


图 6-2 FACTOR 过程的步骤

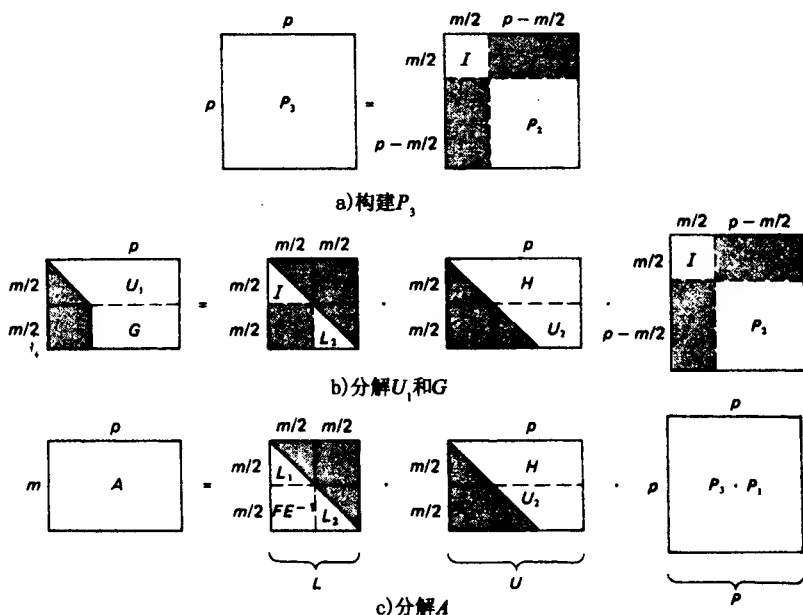


图 6-3 FACTOR 完整过程

⊖ I 可以当作单位下(或上)三角矩阵。

```

procedure FACTOR ( $A, m, p$ );
if  $m = 1$  then
    begin
1.    令  $L = [1]$  (例如:  $L$  是单位  $1 \times 1$  矩阵);
2.    如果可能, 选取  $A$  中具有非零元素的一列  $c$ , 同时, 令  $P$  为  $p \times p$  置换阵, 交换列 1 和列  $c$ ;
        comment 注:  $P = P^{-1}$ ;
3.    令  $U = AP$ ;
4.    return( $L, U, P$ )
    end
else
    begin
5.    采用图 6-2a 所示方法划分  $A$  为  $(m/2) \times p$  的矩阵  $B$  和  $C$ ;
6.    调用 FACTOR( $B, m/2, p$ ) 产生  $L_1, U_1, P_1$ ;
7.    计算  $D = CP_1^{-1}$ ;
        comment 在这里,  $A$  可以写成如图 6-2b 所示的三个矩阵的乘法;
8.    令  $E$  和  $F$  分别为如图 6-2c 所示的  $U_1$  和  $D$  的前  $m/2$  列;
9.    计算  $G = D - FE^{-1}U_1$ ;
        comment 注意: 矩阵  $G$  的前  $m/2$  列全为 0.  $A$  可以写成如图 6-2d 所示的矩阵乘法结果;
10.   令  $G'$  为  $G$  中最右的  $p - m/2$  列;
11.   调用 FACTOR( $G', m/2, p - m/2$ ) 产生  $L_2, U_2, P_2$ ;
12.   令  $P_3$  为  $p \times p$  置换矩阵, 其中  $I_{m/2}$  为其左上角,  $P_2$  为其右下角, 如图 6-3a 所示;
13.   计算  $H = U_1 P_3^{-1}$ ;
        comment 这时由矩阵  $U_1$  和  $G$  组成的矩阵可以表示为图 6-3b 所示。替换图 6-3b 的右边为图 6-2d, 从而  $A$ 
            表示为 5 个矩阵的乘积。其中前两个矩阵为单位下三角矩阵, 第三个为上三角矩阵, 后面两个
            矩阵为置换矩阵。通过求得前两个矩阵乘积, 再求后两个矩阵的乘积得到矩阵  $A$  所期望的分解;
14.   令  $L$  为由  $L_1, O_{m/2}, FE^{-1}$  和  $L_2$  构成的  $m \times m$  矩阵, 如图 6-3c 所示;
15.   令  $U$  为  $m \times p$  矩阵, 且  $H$  为矩阵上半部分,  $O_{m/2}$  和  $U_2$  为下半部分, 如图 6-3c 所示;
16.   令  $P$  为  $P_3 P_1$  乘积;
17.   return( $L, U, P$ );
    end

```

图 6-4 过程 FACTOR

每次递归调用 FACTOR 过程是对 $n \times n$ 矩阵 M 的 $m \times p$ 子阵 A 进行操作。每次调用 m 为 2 的幂, 其中 $m \leq p \leq n$ 。输出结果为如图 6-1 所示的矩阵 L, U 和 P 。□

例 6.3 求下述矩阵的 LUP 分解

$$M = \begin{bmatrix} 0 & 0 & 0 & 1 \\ 0 & 0 & 2 & 0 \\ 0 & 3 & 0 & 0 \\ 4 & 0 & 0 & 0 \end{bmatrix}$$

作为开始, 调用 FACTOR($M, 4, 4$), 则调用如下过程

$$\text{FACTOR}\left(\begin{bmatrix} 0 & 0 & 0 & 1 \\ 0 & 0 & 2 & 0 \end{bmatrix}, 2, 4\right)$$

令该矩阵为 A , 调用 FACTOR($[0 \ 0 \ 0 \ 1], 1, 4$), 返回结果如下:

$$L_1 = [1] \quad U_1 = [1 \ 0 \ 0 \ 0] \text{ 和 } P_1 = \begin{bmatrix} 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 \end{bmatrix}$$

交换第 1 列和第 4 列的位置。

在算法第 7 行, 计算

$$D = CP_1^{-1} = [0 \ 0 \ 2 \ 0] \begin{bmatrix} 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 \end{bmatrix}^{\ominus} = [0 \ 0 \ 2 \ 0]$$

在算法第 8 行, $E = [1]$ 和 $F = [0]$, 所以在第 9 行 $G = D = [0 \ 0 \ 2 \ 0]$ 。在第 10 行, 可得 $G' = [0 \ 2 \ 0]$, 所以在第 11 行得到如下

$$L_2 = [1] \quad U_2 = [2 \ 0 \ 0] \text{ 和 } P_2 = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

下一步, 在算法第 12 行

$$P_3 = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

在第 13 行

$$H = U_2 P_3^{-1} = [1 \ 0 \ 0 \ 0] \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} = [1 \ 0 \ 0 \ 0]$$

由此, $\text{FACTOR}(\begin{bmatrix} 0 & 0 & 0 & 1 \\ 0 & 0 & 2 & 0 \end{bmatrix}, 2, 4)$ 返回

$$L = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, U = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 \end{bmatrix} \text{ 和 } P = P_3 P_1 = \begin{bmatrix} 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \end{bmatrix}$$

现在回到算法第 6 行的 $\text{FACTOR}(M, 4, 4)$, 上述的 L , U 和 P 相应地变成 L_1 , U_1 和 P_1 。则在第 7 行计算下式

$$D = CP_1^{-1} = \begin{bmatrix} 0 & 3 & 0 & 0 \\ 4 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \end{bmatrix} = \begin{bmatrix} 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 4 \end{bmatrix}$$

在算法第 8 行可得

$$E = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix} \text{ 和 } F = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}$$

\ominus 在这个例子中, 所有置换阵的逆恰好是自身。

因此在算法第9行可得

$$G = D = \begin{bmatrix} 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 4 \end{bmatrix}$$

在算法第10行可得

$$G' = \begin{bmatrix} 3 & 0 \\ 0 & 4 \end{bmatrix}$$

读者可自行验证对 $\text{FACTOR}(G', 2, 2)$ 的调用结果为:

$$L_2 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad U_2 = \begin{bmatrix} 3 & 0 \\ 0 & 4 \end{bmatrix} \text{ 和 } P_2 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

由此 P_3 等于 I_4 , 即算法第12行得到的单位矩阵, 同时第13行可得

$$H = U_1 P_3^{-1} = U_1 = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 \end{bmatrix}$$

最后在第14~16行可得

$$L = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}, U = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 \\ 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 4 \end{bmatrix} \text{ 和 } P = \begin{bmatrix} 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \end{bmatrix} \quad \square$$

接下来将对算法6.1的正确性进行分析与证明。

定理6.3 对于任意的非奇异矩阵 A , 算法6.1能计算出 L , U 和 P , 使 $A = LUP$ 。

证明: 对于图6-2和6-3所示的不同情况的分解细节的正确性将留作练习, 只要证明以下两点:

1. 在算法 FACTOR 的第2行总能找到非零列,
2. 在算法第9行 E^{-1} 总是存在。

令 A 为 $m \times n$ 矩阵。对 m 进行归纳, m 是2的幂, 证明如果 A 的秩为 m , FACTOR 将计算出 L , U 和 P , 使 $A = LUP$, L 、 U 和 P 分别为下三角矩阵、上三角矩阵和置换矩阵, 其秩分别为 m , m 和 n 。进一步可得, U 矩阵的前 m 列构成矩阵的秩为 m 。如果 $m=1$, 则 A 一定具有非0元素, 因此归纳假设条件成立。当 $m=2^k$, $k \geq 1$ 时, 因为 A 存在 m 列, 秩为 m , 在算法第5行的 B 和 C 都有 $m/2$ 列, 且各自的秩为 $m/2$ 。由归纳假设, 在算法第6行调用 FACTOR 得到预期的 L_1 , U_1 和 P_1 , 且 U_1 的前 $m/2$ 列秩为 $m/2$ 。因此第9行可得 E^{-1} 存在。

从图6-2d可以看到 A 是三个矩阵的乘积, 其中之一为上面 U_1 下面 G 构成。由于 A 的秩为 m , 所以该矩阵的秩必须是 m 。因此, G 的秩为 $m/2$ 。由于 G 的前 $m/2$ 列为零, 且 G' 是 G 减掉前 $m/2$ 列所得到的矩阵, 因此, G' 的秩也为 $m/2$ 。通过归纳假设, 在算法第11行调用 FACTOR , 得到期望的结果 L_2 , U_2 和 P_2 。通过归纳假设能直接得到结果。 \square

在进行时间分析之前, 观察发现置换矩阵可以采用数组 P 表示, 使得 $P[i] = j$ 当且仅当在第 i 列的第 j 行具有元素1。通过设定 $P_1 P_2[i] = P_1[P_2[i]]$, 两个 $n \times n$ 置换矩阵的乘积可以在 $O(n)$ 时间内计算得到。在这种表示方式中, 置换矩阵的逆同样可以在 $O(n)$ 时间内计算得到。

定理6.4 假定对于每个 n , 可以在 $M(n)$ 时间内计算两个 $n \times n$ 矩阵的乘积, 且对于所有的 m , 存在 $\varepsilon > 0$, 使得 $M(2m) \geq 2^{2+\varepsilon} M(m)$ 。^① 则存在常数 k , 使算法6.1对任意的非奇异矩阵至多需要 $kM(n)$ 时间。

① 直观上看, 这个条件要求 $M(n)$ 在 $n^{2+\varepsilon} \sim n^3$ 之间。因此, 可以说存在某个常数 k , 使 $M(n) = kn^2 \log n$, 在这种情况下, 定理的假设不满足。

证明: 令算法 6.1 运用于 $n \times n$ 矩阵, 且 $T(m)$ 为调用 $\text{FACTOR}(A, m, p)$ 所需要的时间, 其中 A 为 $m \times p$ 矩阵, $m \leq p \leq n$. 由算法 FACTOR 的 1~4 行, 知存在常数 b , 使 $T(1) = bn$. 由于递归的原因, 在算法第 6 行和第 11 行分别需要 $T(m/2)$ 时间, 在算法第 7 行及第 13 行需要做置换矩阵的逆运算, 分别需要 $O(n)$ 时间和任意矩阵乘以置换阵的情况. 仅交换第一个矩阵的相应列. 用数组 P 表示置换阵, 很容易发现第 1 个矩阵的第 $P[i]$ 列将会变成最终矩阵的第 i 列. 该矩阵的乘法结果可以在 $O(mn)$ 时间内得到. 因此在算法第 7 行和第 13 行分别需要时间 $O(mn)$.

在算法第 9 行采用定理 6.2 计算 E^{-1} 需要 $O(M(m/2))$ 时间, 同样的时间用于计算积 FE^{-1} . 由于 U_1 最大为 $(m/2) \times n$, 因此, 积 $(FE^{-1})U_1$ 可在以下时间内计算得到

$$O\left(\frac{n}{m}M\left(\frac{m}{2}\right)\right)$$

注意, m 能整除 n , 因为它们都是 2 的幂, 且 $m \leq n$. 剩下的步骤可以明显地看出最坏情况为 $O(mn)$. 因此, 可以得到以下的递归方程.

$$T(m) \leq \begin{cases} 2T\left(\frac{m}{2}\right) + \frac{cn}{m}M\left(\frac{m}{2}\right) + dmn, & \text{若 } m > 1 \\ bn, & \text{若 } m = 1 \end{cases} \quad (6-5)$$

其中 b, c 和 d 为常数.

通过定理假设及 $M(1) = 1$, 可得 $M(m/2) \geq (m/2)^2$. 因此, 可以结合式 (6-5) 中的第二项和第三项. 存在某个常数 e , 使得

$$T(m) \leq \begin{cases} 2T\left(\frac{m}{2}\right) + \frac{en}{m}M\left(\frac{m}{2}\right), & \text{若 } m > 1 \\ bn, & \text{若 } m = 1 \end{cases} \quad (6-6)$$

从式 (6-6) 可以得到

$$T(m) \leq \frac{en}{4m} \left[4M\left(\frac{m}{2}\right) + 4^2M\left(\frac{m}{2^2}\right) + \cdots + 4^{\log_2 m}M(1) \right] + bnm \leq \frac{en}{4m} \sum_{i=1}^{\log_2 m} 4^i M\left(\frac{m}{2^i}\right) + bnm$$

由定理假设 $4^i M(m/2^i) \leq (1/2^{\epsilon})^i M(m)$, 可得

$$T(m) \leq \frac{en}{4m} M(m) \sum_{i=1}^{\log_2 m} \left(\frac{1}{2^{\epsilon}}\right)^i + bnm$$

由于收敛且 $M(m) \geq m^2$, 存在常数 k , 使 $T(m) \leq (kn/m)M(m)$. 在算法 6.1 中, $n = m$, 因此, $T(n) \leq kM(n)$. \square

推论 给定任意的非奇异矩阵 A , 可以在 $O(n^{2+1})$ 步内找到它的 LUP 分解.

证明: 通过定理 6.1, 6.3 及 6.4 可以得证. \square

6.5 LUP 分解的应用

本节将讲述如何把 LUP 分解运用到矩阵求逆、行列式计算以及联立线性方程组求解等应用中. 可以看出, 上述问题都可规约为两个矩阵的乘法, 因此, 对矩阵乘法渐近时间复杂度的改进都会对这些问题求解的渐近时间复杂度有相应的影响. 反之, 矩阵乘法也可以规约于矩阵的求逆, 因此矩阵乘法与矩阵求逆运算在计算复杂度上是等价的.

定理 6.5 令 $\epsilon > 0$ 且 $a \geq 1$, 令 $M(n)$ 为某个环上的两个矩阵做乘法运算所需的时间, 对于某个 $\epsilon > 0$, $M(2m) \geq 2^{2+\epsilon} M(m)$. 任意非奇异矩阵的逆运算可以在 $O(M(n))$ 时间内完成.

证明: 令 A 是任意 $n \times n$ 的非奇异矩阵, 通过定理 6.3、6.4 可以在 $O(M(n))$ 时间内找到 $A = LUP$. 因此, $A^{-1} = P^{-1}U^{-1}L^{-1}$, P^{-1} 很容易地在 $O(n)$ 步内求出, U^{-1} 和 L^{-1} 存在,

且可以运用定理 6.2 在 $O(M(n))$ 步内得到。积 $P^{-1}U^{-1}L^{-1}$ 可以类似地在 $O(M(n))$ 步内得到。□

推论： $n \times n$ 矩阵的逆可以在 $O(n^{2.81})$ 步内计算得到。

定理 6.6 假设 $M(n)$ 如定理 6.5 所述, A 是 $n \times n$ 矩阵, 则可以在 $O(M(n))$ 步内计算出 $\det(A)$ 。

证明： 运用算法 6.1 找到矩阵 A 的 LUP 分解。假定在算法第 2 行不存在非零列或者在算法第 9 行不存在 E^{-1} 导致算法运行失败, 可得 A 是奇异矩阵且 $\det(A) = 0$ 。否则, 令 $A = LUP$, 则 $\det(A) = \det(L)\det(U)\det(P)$ 。 $\det(L)$ 和 $\det(U)$ 可通过主对角元素乘积获得。由于 L 是单位下三角阵, 可得 $\det(L) = 1$ 。由于 U 是上三角矩阵, 可以在 $O(n)$ 步内得到 $\det(U)$ 。由于 P 是个置换阵, 则 $\det(P) = \pm 1$, 正负号取决于 P 所表示的是奇排列还是偶排列。确定是奇排列还是偶排列通过从排列 $(1, 2, \dots, n)$ 实际进行交换得到, 当前列位置的交换次数, 且至多交换 $n-1$ 次。□

推论： $n \times n$ 矩阵的行列式可以在 $O(n^{2.81})$ 步内得到。

例 6.4 计算例 6.3 中矩阵 M 的行列式。先求出该矩阵的 LUP 分解

$$\begin{bmatrix} 0 & 0 & 0 & 1 \\ 0 & 0 & 2 & 0 \\ 0 & 3 & 0 & 0 \\ 4 & 0 & 0 & 0 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 \\ 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 4 \end{bmatrix} \begin{bmatrix} 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \end{bmatrix}$$

第一个和第二个矩阵的行列式结果通过主对角元素的乘积求得, 分别为 $+1$ 和 $+24$ 。现在仅需要知道最后一个矩阵 P 是奇排列还是偶排列。因为 P 表示排列 $(4, 3, 2, 1)$, 该排列建立两个变换序列为: $(1, 2, 3, 4) \Rightarrow (4, 2, 3, 1) \Rightarrow (4, 3, 2, 1)$, 排列为偶排列, 因此 $\det(P) = +1$, 可得 $\det(M) = +24$ 。□

定理 6.7 令 $M(n)$ 为定理 6.5 所描述, A 为非奇异 $n \times n$ 矩阵, 列向量 b 的长度为 n , x 为未知的列向量 $[x_1, x_2, \dots, x_n]^T$ 。能够在 $O(M(n))$ 步内求解联立线性方程组 $Ax = b$ 。

证明： 运用算法 6.1 写出 $A = LUP$, 则 $LUPx = b$ 的求解分为两步。第一步求解 $Ly = b$, 得到未知向量 y , 第二步采用 $UPx = y$ 求解 x 。上述每个子过程都可以在 $O(n^2)$ 步内使用回代法求出, 即先求解 y_1 , 再用其值代换未知量 y_1 以求解 y_2 , 等等。LUP 分解可以运用定理 6.4 在 $O(M(n))$ 步内完成, 求解 $LUPx = b$ 可以在 $O(n^2)$ 步内完成。□

推论： 具有 n 个未知数的 n 个方程构成的联立方程组可以在 $O(n^{2.81})$ 步内求解。

最后, 将证明矩阵乘法与矩阵求逆的计算复杂度是等价的。

定理 6.8 令 $M(n)$ 和 $I(n)$ 分别为两个 $n \times n$ 矩阵相乘及一个 $n \times n$ 矩阵求逆所需的时间。假定存在某个 $\varepsilon > 0$, 使 $8M(m) \geq M(2m) \geq 2^{2+\varepsilon}M(m)$, $I(n)$ 也类似。则 $M(n)$ 和 $I(n)$ 关于常数因子是等价的。

证明： 定理 6.5 表明, 对某个常数 c_1 , $I(n) \leq c_1M(n)$ 。对于某个 c_2 , 为了建立关系 $M(n) \leq c_2I(n)$, 令 A 和 B 为 $n \times n$ 矩阵, 则

$$\begin{bmatrix} I & A & 0 \\ 0 & I & B \\ 0 & 0 & I \end{bmatrix}^{-1} = \begin{bmatrix} I & -A & AB \\ 0 & I & -B \\ 0 & 0 & I \end{bmatrix}$$

因此, 通过求 $3n \times 3n$ 矩阵的逆来求 AB 的积, 可得 $M(n) \leq I(3n) \leq I(4n) \leq 64I(n)$ 。□

6.6 布尔矩阵的乘法

在 5.9 节中分析了两个布尔矩阵乘法的问题, 其元素选自闭半环 $\{0, 1\}$, 且加法及乘法运

算定义如下

$$\begin{array}{c|cc} + & 0 & 1 \\ \hline 0 & 0 & 1 \\ 1 & 1 & 1 \end{array} \quad \begin{array}{c|cc} \cdot & 0 & 1 \\ \hline 0 & 0 & 0 \\ 1 & 0 & 1 \end{array}$$

可以证明两个布尔矩阵做乘法运算在时间复杂度上等价于求图的传递闭包,但遗憾的是,闭半环 $\{0, 1\}$ 不是环,因此,Strassen 矩阵乘法算法及至本章为止所提到的结果在此不能直接运用于布尔矩阵的乘法。

很明显,通常矩阵乘法需要 $O(n^3)$ 步。[⊖]但至少还存在两种时间复杂度低于 $O(n^3)$ 步的布尔矩阵乘法。第一种方法的时间渐近复杂度表现更为优良,但第二种在 n 适中的时候更切合实际。下面定理中将展示第一种方法。

定理 6.9 可以在 $O_A(n^{2.81})$ 步内获得两个 $n \times n$ 布尔矩阵 A 和 B 的乘积。

证明: 模 $n+1$ 的整数形成环 Z_{n+1} , 运用 Strassen 矩阵乘法算法在环 Z_{n+1} 中计算 A 和 B 的乘积。令 C 为该乘积, D 为该乘积对应的布尔矩阵。可以看出: 如果 $D[i, j] = 0$, 则 $C[i, j] = 0$; 如果 $D[i, j] = 1$, 则 $1 \leq C[i, j] \leq n$ 。因此, D 可以容易从 C 得到。□

推论 1 假如两个 k 位整数乘法需要 $m(k)$ 次按位操作, 则布尔矩阵乘法可以在 $O_B(n^{2.81} m(\log n))$ 步内完成。

证明: 由于所有的运算都基于 Z_{n+1} , 因此至多需要 $\lfloor \log n \rfloor + 1$ 位来表示数值。因此这样的两个整数相乘至多需要 $O_B(m(\log n))$ 时间, 且加法或者减法至多需要 $O_B(\log n)$ 时间, 因此不可能再大, 从而得证。□

在第 7 章中, 将讨论一种 $m(k)$ 为 $O_B(k \log k \log \log k)$ 的乘法算法, 运用这个值得到如下推论。

推论 2 布尔矩阵乘法至多需要 $O_B(n^{2.81} \log n \log \log n \log \log \log n)$ 步。

第二种方法常称为“4 俄罗斯人”(Four Russians)算法, 以发明者人数和国籍命名, 从某种意义上比定理 6.9 中的算法更为“实际”。而且, 该算法更容易进行位向量运算, 算法 6.9 没有这个性质。

假定计算 A 和 B 的乘积, 它们是两个 $n \times n$ 布尔矩阵。为了方便, 假定 $\log n$ 能够整除 n , 则可以把 A 分成 $n \times (\log n)$ 个子阵, 同时 B 可以分成 $(\log n) \times n$ 个子阵, 如图 6-5 所示。计算可写成:

$$AB = \sum_{i=1}^{n/\log n} A_i B_i$$

注意, 积 $A_i B_i$ 本身为 $n \times n$ 矩阵, 假设可以在 $O(n^2)$ 步内计算出 $A_i B_i$ 的积, 那么, 可以在 $O(n^3/\log n)$ 步内计算出 AB , 因为具有 $n/\log n$ 个这样的积。

现在重点讨论积 $A_i B_i$ 的计算。可以在 $O(n^2 \log n)$ 步内计算出来。为了计算 $A_i B_i$, 可以对 A_i 中的每行 a_j 计算 $a_j B_i$ 。为了计算 $a_j B_i$, 可以找到 a_j 对应在 k 列有 1 的 B_i 的每个 k 行。然后将 B_i 那些行加起来, 形成长为 n 的行的位向量。

例如, 假定

⊖ 如果积的所有元素几乎都是 1 的话, 则可以在 $O(n^2)$ 期望时间内计算两个布尔矩阵的乘法, 计算 $\sum_{k=1}^n a_{ik} b_{kj}$, 直到有一项为 1 为止。由此可知和为 1。举例来说, 如果每个元素 a_{ik} 或 b_{kj} 为 1 的概率分别为 p , 则可以找到此类项的期望值至多为 $1/p^2$, 与 n 无关。

$$A_i = \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 1 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \\ 1 & 1 & 0 \\ 0 & 0 & 0 \\ 0 & 1 & 1 \end{bmatrix} \text{ 和 } B_i = \begin{bmatrix} 0 & 1 & 0 & 1 & 1 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 & 1 & 0 & 0 \\ 1 & 1 & 0 & 1 & 0 & 0 & 0 & 0 \end{bmatrix}$$

可得

$$C_i = A_i B_i = \begin{bmatrix} 1 & 1 & 0 & 1 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 1 & 1 & 0 & 0 & 1 \\ 1 & 1 & 0 & 1 & 1 & 1 & 0 & 1 \\ 0 & 1 & 0 & 1 & 1 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 1 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 1 & 0 & 1 & 0 & 0 \end{bmatrix}$$

C_i 的第一行仅为 B_i 的第三行, 因为 A_i 的第一行仅在第三列有一个 1。 C_i 的第二行为 B_i 的第一行和第三行的和, 因为 A_i 的第二行在第 1 列和第 3 列是 1, 等等。

为了计算 $a_i B_i$, 对于 A_i 的每行 a_i 需要 $O(n \log n)$ 时间, 由于 A_i 有 n 行, 因此计算 $A_i B_i$ 所需的总时间为 $O(n^2 \log n)$ 。

为了更快地计算积 $A_i B_i$, 可以观察 A_i 的各行具有 $\log n$ 个元素, 且 A_i 中各元素要么为 0 要么为 1。因此, 至多只有 $2^{\log n} = n$ 个不同的行。

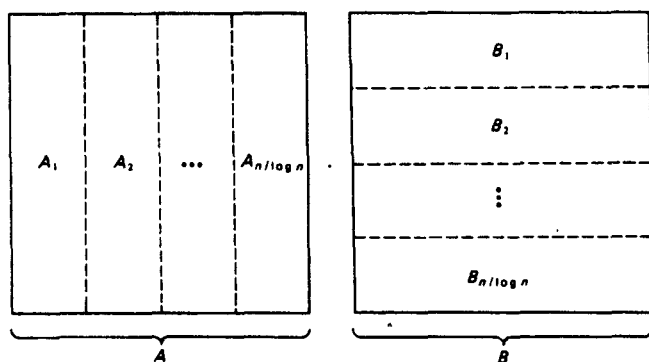


图 6-5 布尔矩阵 A 和 B 的划分

因此, B_i 各行的可能和仅有 n 个不同的值。可以预先把 B_i 各行的可能和求出来, 不计算 $a_i B_i$, 而仅用 a_i 的表简单索引查询相应的结果。

这种方法仅需要 $O(n^2)$ 时间计算 $A_i B_i$ 。原因如下。 B_i 行的子集有空集, 单元素的集合, 或者是单元素集合与一个小集合的并集。可以通过选择恰当的计算顺序, 把一行加到前面已经得到结果的基础上。因此, 可以在 $O(n^2)$ 步内计算 B_i 中行的所有 n 个和。在此和的基础上, 把它们放在数组中, 针对 A_i n 行的每个不同行选择相应的和。

算法 6.2 4 俄罗斯人布尔矩阵乘法算法。

输入：两个 $n \times n$ 的布尔矩阵 A 和 B 。

输出：积 $C = AB$

方法：令 $m = \lfloor \log n \rfloor$ ，把 A 划分成 $A_1, A_2, \dots, A_{\lceil n/m \rceil}$ ，其中 $A_i (1 \leq i < \lceil n/m \rceil)$ 表示矩阵 A 由第 $m(i-1)+1$ 列到第 mi 列构成的矩阵，且 $A_{\lceil n/m \rceil}$ 由剩下的列构成，通过补充 0 列达到必需的 m 列。划分矩阵 B 为 $B_1, B_2, \dots, B_{\lceil n/m \rceil}$ ，其中 $B_i (1 \leq i < \lceil n/m \rceil)$ 表示矩阵 B 从第 $m(i-1)+1$ 行到第 mi 行构成的矩阵， $B_{\lceil n/m \rceil}$ 由剩下的行所组成。如果不够 m 行，通过补 0 行实现。所述情况如图 6-5 所示。计算过程如图 6-6 所示。采用 $\text{NUM}(v)$ 表示 0 和 1 的向量 v 的逆。例如， $\text{NUM}([0, 1, 1]) = 6$ 。□

```

begin
1.   for i ← 1 until ⌈n/m⌉ do
       begin
           comment 计算  $B_i$  中各行的和，其中各行为  $b_i^{(0)}, \dots, b_i^{(m)}$ ;
2.       ROWSUM[0] ← [0, 0, ..., 0]; 0 的个数为  $n$ 
3.       for j ← 1 until  $2^m - 1$  do
           begin
4.               令  $k$  为符合不等式  $2^k \leq j < 2^{k+1}$  的值;
5.               ROWSUM[j] ← ROWSUM[j -  $2^k$ ] +  $b_{k+1}^{(j)}$ ①
           end;
6.       令  $C_i$  为矩阵且其第  $j$  行 ( $1 \leq j \leq n$ ) 为 ROWSUM[ $\text{NUM}(a_j)$ ],
           其中  $a_j$  为  $A_i$  的第  $j$  行
       end;
7.   let  $C$  be  $\sum_{i=1}^{\lceil n/m \rceil} C_i$ 
end

```

① 当然，这里的求和是按位布尔操作。

图 6-6 4 俄罗斯人算法

定理 6.10 计算 $C = AB$ 的算法 6.2 可在 $O(n^3/\log n)$ 步内完成。

证明：对 j 进行归纳可得算法第 2~5 行的 $\text{ROWSUM}[j]$ 是 j 的二进制表示方法中具有 1 的位置上对应 B_i 中 b_k 行按位布尔求和的值，该位置是从右往左进行计算的第 k 个位置。在算法第 6 行得到 $C_i = A_i B_i$ ，然后在第 7 行得到 $C = AB$ 。

在算法时间复杂度方面，首先考虑 3~5 行的循环。第 5 行的赋值语句明显需要 $O(n)$ 步。在第 4 行计算 k 需要 $O(m)$ ，比 $O(n)$ 小，所以第 4~5 行整体循环复杂度为 $O(n)$ 。这个循环要重复 $2^m - 1$ 次，所以循环复杂度是 $O(n2^m)$ 。由于 $m \leq \log n$ ，则第 3~5 行的循环复杂度为 $O(n^2)$ 。□

在第 6 行，计算 $\text{NUM}(a_j)$ 的时间复杂度是 $O(m)$ ，复制向量 $\text{ROWSUM}[\text{NUM}(a_j)]$ 的复杂度为 $O(n)$ ，因此，第 6 行的复杂度是 $O(n^2)$ 。由于 $\lceil n/m \rceil \leq 2n/\log n$ ，且第 1~6 行的循环执行 $\lceil n/m \rceil$ 次，则时间复杂度是 $O(n^3/\log n)$ 。同样，第 7 步至多需要 $2n/\log n$ 来完成两个 $n \times n$ 矩阵的求和，总共需要 $O(n^3/\log n)$ 步。因此，整个算法的时间复杂度为 $O(n^3/\log n)$ 步。□

更为有趣的事情是若具有针对位串进行操作的逻辑和算术指令，则算法 6.2 只需要 $O_{\text{bv}}(n^2/\log n)$ 的位向量步运算。

定理 6.11 算法 6.2 可以在 $O_{\text{bv}}(n^2/\log n)$ 的位向量步内完成运算。

证明：设置累加器用来确定何时增加 k 。初始累加器的值为 1， k 的值为 0。在 j 每次增加时，累加器的值减 1，除非值为 1，此时累加器的值设置为 j 的新值，且 k 的值增加 1。

在图 6-6 中算法第 2 行和第 5 行需要常数步。因此，第 3~5 行的循环是 $O_{\text{bv}}(n)$ 。在算法第

6行建立 C_i , 注意, 由于在 RAM 中可用整数表示位向量, 因此计算 $\text{NUM}(a_i)$ 不需要时间。[⊖] 因此, C_i 的每行可以在常数位向量步内找到, 第6行需要 $O_{BV}(n)$ 。因此第1~6行的循环复杂度是 $O_{BV}(n^2/\log n)$, 且第7行的复杂度一样。□

习题

- 6.1 证明整数模 n 运算形成环。即 Z_n 是环 $(\{0, 1, \dots, n-1\}, +, \cdot, 0, 1)$, 其中 $a+b$ 和 $a \cdot b$ 为模 n 的普通加法和乘法运算。
- 6.2 证明元素选自某个环 R 的 $n \times n$ 矩阵的集合 M_n 形成环。
- 6.3 给出例子证明矩阵的乘法不符合交换律, 即使其元素选自某个环, 且其中乘法运算符合交换律。
- 6.4 用 Strassen 算法计算乘法:

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix}$$

- 6.5 Strassen 算法的另一个版本采用以下恒等式帮助计算两个 2×2 矩阵的积。

$$\begin{aligned} s_1 &= a_{21} + a_{22} & m_1 &= s_2 s_6 & t_1 &= m_1 + m_2 \\ s_2 &= s_1 - a_{11} & m_2 &= a_{11} b_{11} & t_2 &= t_1 + m_4 \\ s_3 &= a_{11} - a_{21} & m_3 &= a_{12} b_{21} \\ s_4 &= a_{12} - s_2 & m_4 &= s_3 s_7 \\ s_5 &= b_{12} - b_{11} & m_5 &= s_1 s_5 \\ s_6 &= b_{22} - s_5 & m_6 &= s_4 b_{22} \\ s_7 &= b_{22} - b_{12} & m_7 &= a_{22} s_8 \\ s_8 &= s_6 - b_{21} \end{aligned}$$

结果的矩阵元素为:

$$\begin{aligned} c_{11} &= m_2 + m_3 \\ c_{12} &= t_1 + m_5 + m_6 \\ c_{21} &= t_2 - m_7 \\ c_{22} &= t_2 + m_5 \end{aligned}$$

证明这些元素计算了等公式(6-1), 且仅需要7次乘法和15次加法运算。

- 6.6 证明对于 $n \times n$ 矩阵 A , B 和 C , 以下公式成立。
 - a) $AB = I$, $AC = I$ 可得 $B = C$
 - b) $A^{-1}A = I$
 - c) $(AB)^{-1} = B^{-1}A^{-1}$
 - d) $(A^{-1})^{-1} = A$
 - e) $\det(AB) = \det(A)\det(B)$
- 6.7 定理6.2表明非奇异上三角矩阵的逆对于某个 C , 可以在 $cn^{2.81}$ 内完成算术运算。假定 n 为2的幂且采用 Strassen 算法进行矩阵乘法, 找出常数 c 。
- 6.8 用数组 P 表示置换矩阵, 当且仅当第 i 列的第 j 行为1, $P[i] = j$ 。令 P_1 和 P_2 为两个 $n \times n$ 置换矩阵的表示数组。

⊖ 表示方法的详细说明: $\text{NUM}(a_i)$ 是表示 a_i 的反序的整数, 取自 B_i 的“第 j 行”是从底部开始的第 j 行, 而不是前面所提及的从顶部开始的。

a) 证明 $P_1 P_2[i] = P_1[P_2[i]]$ 。

b) 给出一个 $O(n)$ 算法计算 P_1^{-1} 。

c) 改变表示方式, 当且仅当第 i 行的第 j 列为 1 时, $P[i] = j$, 给出正确的 $P_1 P_2$ 公式, 同时给出计算 P_1^{-1} 的算法。

6.9 用算法 6.1 计算以下矩阵的 LUP 分解。

$$M = \begin{bmatrix} 0 & 0 & 1 & 2 \\ 0 & 0 & 3 & 0 \\ 1 & -1 & 0 & 1 \\ 2 & 0 & -1 & 3 \end{bmatrix}$$

6.10 已经证明 LUP 分解、矩阵求逆、行列式的计算以及线性方程组的求解都是 $O_A(n^{2.81})$ 。假定采用 Strassen 算法做矩阵乘法, n 是 2 的幂, 采用算法 6.1 及定理 6.4 ~ 6.7, 分别找到最合适的常数因子。

6.11 用本章讲述的方法求习题 6.9 中矩阵 M 的逆及行列式。

6.12 用 LUP 分解求解下述联立线性方程组:

$$x_3 + 2x_4 = 7$$

$$3x_3 = 9$$

$$x_1 - x_2 + x_4 = 3$$

$$2x_1 - x_3 + 3x_4 = 10$$

6.13 证明任意的排列不是奇排列就是偶排列, 但不可能同时是两种排列。

6.14 分别采用定理 6.9 中的算法及 4 俄罗斯人算法计算下面布尔矩阵的乘积。

$$\left[\begin{array}{cccc|cccc} 1 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 & 1 & 1 & 0 & 0 \\ 1 & 0 & 0 & 1 & 1 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & 1 \end{array} \right]$$

6.15 通过证明图 6-2 及图 6-3 间的有效关系来完成定理 6.3 的证明。

** 6.16 考虑模 2 运算的整数域 $^{\odot}F_2$, 寻找一种基于 F_2 的 $n \times n$ 矩阵乘法算法, 渐近复杂度为 $n^{2.81}/(\log n)^{0.4}$ [提示: 把矩阵划分成大小为 $\sqrt{\log n} \times \sqrt{\log n}$ 的矩阵块]。

6.17 估算 n 的值, 使之当超过该数值时, $n^{2.81} < n^3/\log n$ 。

* 6.18 令 $L(n)$ 为两个 $n \times n$ 下三角矩阵相乘所需的时间, $T(n)$ 为任意两个矩阵相乘所需的时间。证明存在常数 c , 使 $T(n) \leq cL(n)$ 。

6.19 证明上(下)三角矩阵的逆还是上(下)三角矩阵。

* 6.20 令 $I(n)$ 和 $U(n)$ 分别为求任意 $n \times n$ 矩阵及 $n \times n$ 上三角矩阵的逆所需的时间步数。证明存在常数 c , 对于所有 n , 使 $I(n) \leq cU(n)$ 。

** 6.21 为了计算矩阵积 $C = AB$, 可以先计算积 $D = (PAQ)(Q^{-1}BR)$, 然后计算 $C = P^{-1}DR^{-1}$ 。如果 P , Q 和 R 为某些矩阵, 比如置换矩阵, 乘法 PAQ 、 $Q^{-1}BR$ 及 $P^{-1}DR^{-1}$ 不需要环元素的乘法运算。采用这种思想寻找另外一种 2×2 矩阵的乘法算法, 只需要 7 次乘法操作。

6.22 证明当 LU 存在时, 非奇异矩阵 A 的 LU 分解是唯一的 [提示: 假定 $A = L_1 U_1 = L_2 U_2$, 证明 $L_2^{-1} L_1 = U_2 U_1^{-1} = I$]。

6.23 证明: 如果 A 是非奇异矩阵, 且 A 的任意主子阵是非奇异的, 则 A 存在 LUP 分解。

6.24 奇异矩阵具有 LUP 分解吗?

- ** 6.25** 假定 A 是实数的 $n \times n$ 矩阵, 如果对于任意的非零向量 x , 使得 $x^T A x > 0$, 则 A 为正定矩阵。
- 证明引理 6.5 可以用于求取任意非奇异的对称正定矩阵的逆。
 - 证明 AA^T 总是正定和对称的。
 - 采用 a) 和 b) 中的结论, 给出一个 $O(M(n))$ 算法用于求任意实数的非奇异矩阵的逆。
 - (c) 中的算法对模 2 的整数域有效吗?
- * 6.26** $n \times n$ 特普利茨 (Toeplitz) 矩阵 A 具有如下性质
- $$A[i, j] = A[i-1, j-1], \quad 2 \leq i, j \leq n$$
- 寻找特普利茨矩阵的表示方式, 使两个 $n \times n$ 特普利茨矩阵相加可以在 $O(n)$ 次操作内完成。
 - 寻找 $n \times n$ 的特普利茨矩阵和列向量相乘的分治算法。需要多少次算术操作? [提示: 采用分治算法可以获得 $O(n^{1.59})$ 的结果。第 7 章将采用一种新技术来优化这个结果。]
 - 寻找有效的渐近算法用于两个 $n \times n$ 特普利茨矩阵相乘。该算法需要多少次算术操作? 注意特普利茨矩阵相乘的结果不一定是特普利茨矩阵。

研究性问题

- 6.27** 一个固有的问题是直接提高 Strassen 算法的效率。Hopcroft 和 Kerr[1971] 已经证明基于任意环的 2×2 矩阵乘法需要 7 次乘法操作。但递归算法可能基于另外更小规模的矩阵。例如, 如果 3×3 矩阵乘法能够在 21 次乘法操作内完成, 或者 4×4 矩阵在 48 次操作内完成, Strassen 算法的效率就可以渐近提高。
- 6.28** 最短路径问题能够在低于 $O(n^3)$ 步内完成吗? Strassen 算法并不适合由包括 $+\infty$ 的非负实数构成的闭半环, 但在做布尔矩阵运算时可以把闭半环嵌入到一个环中。

文献和注释

Strassen 算法取自 Strassen[1969]。Winograd[1973] 把加法次数降为 15 次, 该方法优化了常数因子, 但在复杂度的阶上没有变化(见习题 6.5)。

Strassen[1969] 同样给出了时间复杂度为 $O(n^{2.81})$ 的矩阵求逆、行列式计算和联立线性方程组求解的算法, 在此假定所有矩阵均为非奇异矩阵。Bunch 和 Hopcroft[1974] 给出在 $O(n^{2.81})$ 步内进行 LUP 分解的算法, 前提是初始矩阵为非奇异矩阵。A. Schönhage 独自证明了基于有序域非奇异矩阵的求逆运算可以在 $O(n^{2.81})$ 步内完成(习题 6.25)。

矩阵乘法不难于矩阵求逆的结论由 Winograd[1970c] 得到。Fischer 和 Meyer[1971] 给出布尔矩阵乘法的 $O_A(n^{2.81})$ 算法。“4 俄罗斯人”算法由 Arlazarov, Dinic, Kronrod 和 Faradzev[1970] 给出。

更多关于矩阵的数学方法, 可以查阅 Hohn[1958]。一些代数概念, 如环理论可以在 MacLane 和 Birkhoff[1967] 中找到。习题 6.13 的解决方案可以在 Birkhoff 和 Bartee[1970] 中找到。习题 6.16 由 J. Hopcroft 提出。

第7章 快速傅里叶变换及其应用

傅里叶变换已广泛应用于科学和工程领域,所以,计算傅里叶变换的有效算法很重要。另外,傅里叶变换很普遍,很有必要设计一个高效的傅里叶变换算法。在许多应用中,习惯于把一个问题转化为另外一个比较简单的问题,如两个多项式乘积的计算。首先对多项式系数向量进行一次线性变换,再对系数向量进行一次简单的卷积操作,最后再进行一次逆变换便可得到多项式的期望乘积。这种线性变换就是离散傅里叶变换。

本章主要研究傅里叶变换、傅里叶的逆变换以及它们在计算卷积和各种类型乘积中的作用。本章将设计一个高效的算法,即快速傅里叶变换(FFT),通过除法计算多项式的值,同时利用多项式在单位根上计值的特点。

还将证明卷积理论。把卷积解释为单位根的多项式计算、样本值的乘法和多项式插值。我们将利用快速傅里叶变换设计一个高效的卷积算法,然后把高效的卷积算法用到多项式的符号乘法和整数的乘法运算中。所产生的整数乘法算法,所谓 Schönhage - Strassen 算法是已知的渐近最快的两个整数相乘的方法。

7.1 离散傅里叶变换及其逆变换

傅里叶变换常定义在复数域上。把傅里叶变换定义在任意交换环(commutative ring) $(R, +, \cdot, 0, 1)^\ominus$, 其原因在后面将会明了。 R 中的元素 ω 定义如下:

1. $\omega \neq 1$,
2. $\omega^n = 1$, 并且
3. $\sum_{j=0}^{n-1} \omega^{jp} = 0 (1 \leq p < n)$

元素 ω 称为主 n 次单位根(principal n th root of unity)。元素 $\omega^0, \omega^1, \dots, \omega^{n-1}$ 分别是 n 次单位根(n th roots of unity)。

例如: $e^{2\pi i/n}$ (其中 $i = \sqrt{-1}$) 是复数环中的主 n 次单位复根。

假设 $a = [a_0, a_1, \dots, a_{n-1}]^T$ 是一个长度为 n 的(列)向量,各个分量均是集合 R 中的元素。假设整数 n 在 R 上存在乘法的逆[⊖], 并且 R 存在主 n 次单位复根 ω 。假设 A 是一个 $n \times n$ 矩阵,其中 $A[i, j] = \omega^{ij} (0 \leq i, j < n)$ 。向量 $F(a) = Aa = \sum_{k=0}^{n-1} a_k \omega^{ik}$ 称为 a 的离散傅里叶变换,其第 i 个分量是 $b_i (0 \leq i < n)$ 。若矩阵 A 是非奇异的,则 A^{-1} 存在,且 A^{-1} 有着引理 7.1 定义的简单形式。

引理 7.1 假设 R 是一个交换环,且存在主 n 次单位根 ω , 其中 n 在 R 中存在乘法的逆。假设 A 是一个 $n \times n$ 矩阵,其第 ij 个元素是 $\omega^{ij} (0 \leq i, j < n)$ 。那么 A^{-1} 存在,且第 ij 个元素是 $(1/n)\omega^{-ij}$ 。

证明: 假设当 $i=j$ 时 $\delta_{ij} = 1$, 其他情况下为 0。很显然, A^{-1} 的定义如上面所述, $A \cdot A^{-1} = I_n$, 也就是 $A \cdot A^{-1}$ 的第 ij 个元素是

⊖ 交换环是一个乘法(和加法)符合交换律的环。

⊖ 整数出现在任何环中,即使是有限环。把 n 写成 $1+1+\dots+1$ (n 次), 其中 1 是乘法单位元。也就是 $\omega\omega^{-1}=1$, 说明 ω 存在逆元,这样就可以讨论 ω 的负幂。

$$\frac{1}{n} \sum_{k=0}^{n-1} \omega^{ik} \omega^{-kj} = \delta_{ij} \quad \text{对于 } 0 \leq i, j < n \quad (7-1)$$

如果 $i=j$, 公式(7-1)左边化简为

$$\frac{1}{n} \sum_{k=0}^{n-1} \omega^0 = 1$$

如果 $i \neq j$, 假设 $q=i-j$, 公式(7-1)左边可以化简为

$$\frac{1}{n} \sum_{k=0}^{n-1} \omega^{qk} \quad (-n < q < n, q \neq 0)$$

如果 $q > 0$, 则

$$\frac{1}{n} \sum_{k=0}^{n-1} \omega^{qk} = 0$$

因为 ω 是主 n 次单位根。如果 $q < 0$, 那么乘以 $\omega^{-q(n-1)}$, 重新调整和式中各项的顺序, 用 $-q$ 替换 q 可以得到

$$\frac{1}{n} \sum_{k=0}^{n-1} \omega^{qk} \quad (0 < q < n)$$

由于 ω 是主 n 次单位根, 所以它的值为 0。方程(7-1)马上可以得证。 \square

向量 $F^{-1}(\mathbf{a}) = A^{-1}\mathbf{a}$, 其第 i ($0 \leq i < n$) 个分量为

$$\frac{1}{n} \sum_{k=0}^{n-1} a_k \omega^{-ik}$$

是 \mathbf{a} 的离散傅里叶变换的逆。很显然, \mathbf{a} 变换的逆变换就是 \mathbf{a} 本身, 也就是 $F^{-1}F(\mathbf{a}) = \mathbf{a}$ 。

在傅里叶变换和多项式计值/插值之间存在紧密的关系, 假设

$$p(x) = \sum_{i=0}^{n-1} a_i x^i$$

是 $(n-1)$ 次多项式。这个多项式可以用两种方式唯一地表示, 一种表示为系数列表 a_0, a_1, \dots, a_{n-1} , 或者在 n 个确定点 x_0, x_1, \dots, x_{n-1} 列出它的值。找到表示多项式系数值 x_0, x_1, \dots, x_{n-1} 的过程称为插值。

计算向量 $[a_0, a_1, \dots, a_{n-1}]^T$ 的傅里叶变换等同于把多项式 $\sum_{i=0}^{n-1} a_i x^i$ 的系数表示转换到在点 $\omega^0, \omega^1, \dots, \omega^{n-1}$ 上的值表示。类似地, 傅里叶变换的逆等同于在多项式的第 n 次单位根上进行插值操作。

可以定义在点集合, 而不是在第 n 次单位根上的变换。例如, 可以使用整数 $1, 2, \dots, n$ 。然而, 通过对 ω 的幂选择, 计值和插值操作会变得非常直接。在第 8 章中, 傅里叶变换用到任意点上对多项式进行计值和插值。

傅里叶变换的一个主要应用是计算两个向量的卷积操作。假设

$$\mathbf{a} = [a_0, a_1, \dots, a_{n-1}]^T \text{ 和 } \mathbf{b} = [b_0, b_1, \dots, b_{n-1}]^T$$

是两个列向量。 \mathbf{a} 和 \mathbf{b} 向量的卷积操作表示为 $\mathbf{a} \otimes \mathbf{b}$, 结果是向量 $\mathbf{c} = [c_0, c_1, \dots, c_{2n-1}]^T$, 其中 $c_i = \sum_{j=0}^{n-1} a_j b_{i-j}$ (如果 $k < 0$ 或者 $k \geq n$, 假设 $a_k = b_k = 0$)。则

$$c_0 = a_0 b_0, c_1 = a_0 b_1 + a_1 b_0, c_2 = a_0 b_2 + a_1 b_1 + a_2 b_0$$

等等。注, $c_{2n-1} = 0$, 在这里保留它的目的是为了对称。

为了说明卷积操作, 再回顾一下多项式的系数表示法。两个 $(n-1)$ 次多项式的乘积

$$p(x) = \sum_{i=0}^{n-1} a_i x^i \text{ 和 } q(x) = \sum_{j=0}^{n-1} b_j x^j$$

是 $(2n-2)$ 次多项式

$$p(x)q(x) = \sum_{i=0}^{2n-2} \left[\sum_{j=0}^i a_j b_{i-j} \right] x^i$$

从上面可以看出, 如果忽略 c_{2n-1} ($c_{2n-1}=0$), 那么多项式乘积的系数正好是两个原始多项式的系数向量 $[a_0, a_1, \dots, a_n]^T$ 和 $[b_0, b_1, \dots, b_n]^T$ 卷积的分量,

如果用系数来表示两个 $(n-1)$ 次多项式, 那么可以通过计算两个系数向量卷积的方式来计算其乘积的系数。另一方面, 如果用第 n 次单位根上的值来表示 $p(x)$ 和 $q(x)$, 那么, 可以简单地把对应根上的对值相乘来计算它们乘积的表示值。这说明两个向量 \mathbf{a} 和 \mathbf{b} 的卷积是两个向量变换的分量乘积的逆变换。形式化表示为: $\mathbf{a} \otimes \mathbf{b} = F^{-1}(F(\mathbf{a}) \cdot F(\mathbf{b}))$ 。也就是说, 卷积运算可以先使用傅里叶变换, 然后计算它们对应的乘积, 最后再进行一次逆变换得到。唯一存在的问题在于, 两个 $(n-1)$ 次多项式的乘积通常是一个 $(2n-2)$ 次多项式, 要表示它需要 $2n-2$ 个不同点的值。这个技术在下面的定理中得到了解决, 考虑 $p(x)$ 和 $q(x)$ 是 $(2n-1)$ 次多项式的情形, 其中 x 的最高次 $(n-1)$ 的系数是 0 [即, 把 $(n-1)$ 次多项式当作 $(2n-1)$ 次多项式来处理]。

定理 7.1 卷积定理 (Convolution Theorem) 假设

$$\mathbf{a} = [a_0, a_1, \dots, a_{n-1}, 0, \dots, 0]^T$$

和

$$\mathbf{b} = [b_0, b_1, \dots, b_{n-1}, 0, \dots, 0]^T$$

是长度为 $2n$ 的列向量。假设

$$F(\mathbf{a}) = [a'_0, a'_1, \dots, a'_{2n-1}]^T$$

和

$$F(\mathbf{b}) = [b'_0, b'_1, \dots, b'_{2n-1}]^T$$

是它们的傅里叶变换, 那么 $\mathbf{a} \otimes \mathbf{b} = F^{-1}(F(\mathbf{a}) \cdot F(\mathbf{b}))$ 。

证明: 因为 $a_i = b_i = 0$ ($n \leq i < 2n$), 当 $0 \leq l < 2n$ 时

$$a'_l = \sum_{j=0}^{n-1} a_j \omega^{jl} \text{ 和 } b'_l = \sum_{k=0}^{n-1} b_k \omega^{kl}$$

那么

$$a'_l b'_l = \sum_{j=0}^{n-1} \sum_{k=0}^{n-1} a_j b_k \omega^{l(j+k)} \quad (7-2)$$

假设 $\mathbf{a} \otimes \mathbf{b} = [c_0, c_1, \dots, c_{2n-1}]^T$ 和 $F(\mathbf{a} \otimes \mathbf{b}) = [c'_0, c'_1, \dots, c'_{2n-1}]^T$ 。

因为 $c_p = \sum_{j=0}^{2n-1} a_j b_{p-j}$, 所以有

$$c'_l = \sum_{p=0}^{2n-1} \sum_{j=0}^{2n-1-j} a_j b_{p-j} \omega^{lp} \quad (7-3)$$

交换式(7-3)中加法顺序并且用 k 替换 $p-j$ 得到

$$c'_l = \sum_{j=0}^{2n-1} \sum_{k=-j}^{2n-1-j} a_j b_k \omega^{l(j+k)} \quad (7-4)$$

因为 $b_k = 0$ ($k < 0$), 可以把内层加法下限增大到 $k=0$ 。类似地, 因为 $a_j = 0$ ($j \geq n$), 把外层连加的上限缩小到 $n-1$ 。无论 j 的值是多少, 内部连加的上限至少为 n 。由于 $b_k = 0$ ($k \geq n$), 可以把上限替换为 $n-1$ 。通过这些修改, 公式(7-4)就和式(7-2)相同了, 所以有 $c'_l = a'_l b'_l$ 。我们可以得到 $F(\mathbf{a} \otimes \mathbf{b}) = F(\mathbf{a}) \cdot F(\mathbf{b})$, 即 $\mathbf{a} \otimes \mathbf{b} = F^{-1}(F(\mathbf{a}) \cdot F(\mathbf{b}))$ 成立。□

两个长度为 n 的向量的卷积是长度为 $2n$ 的向量, 这需要卷积定理中对向量 \mathbf{a} 和向量 \mathbf{b} “填充” n 个 0。为了避免这项“填充”, 可使用“包”卷积。

定义 假设 $\mathbf{a} = [a_0, a_1, \dots, a_{n-1}]^T$ 和 $\mathbf{b} = [b_0, b_1, \dots, b_{n-1}]^T$ 是两个长度为 n 的向量。 \mathbf{a}

和 \mathbf{b} 的正包卷积 (positive wrapped convolution) 是向量 $\mathbf{c} = [c_0, c_1, \dots, c_{n-1}]^T$, 其中

$$c_i = \sum_{j=0}^i a_j b_{i-j} + \sum_{j=i+1}^{n-1} a_j b_{n+i-j}$$

\mathbf{a} 和 \mathbf{b} 的负包卷积 (negative wrapped convolution) 是向量 $\mathbf{d} = [d_0, d_1, \dots, d_{n-1}]^T$, 其中

$$d_i = \sum_{j=0}^i a_j b_{i-j} - \sum_{j=i+1}^{n-1} a_j b_{n+i-j}$$

稍后, 7.5 节中的 Schönhage-Strassen 算法, 一个快速整数乘法算法用到该卷积。而由目前观察所知, 可以在第 n 个单位根上对两个 $n-1$ 次多项式计值, 并且把对应根上的值相乘。这给出了 n 个值, 通过这些值可以对唯一的 $n-1$ 次多项式进行插值。这个唯一多项式的系数向量正好是两个原始多项式系数向量的正包卷积。

定理 7.2 假设 $\mathbf{a} = [a_0, a_1, \dots, a_{n-1}]^T$ 和 $\mathbf{b} = [b_0, b_1, \dots, b_{n-1}]^T$ 是两个长度为 n 的向量。 ω 是主 n 次单位根, 并且 $\psi^2 = \omega$ 。假设 n 有一个乘法的逆。

1. \mathbf{a} 和 \mathbf{b} 的正包卷积可以由 $F^{-1}(F(\mathbf{a}) \cdot F(\mathbf{b}))$ 计算。

2. 假设 $\mathbf{d} = [d_0, d_1, \dots, d_{n-1}]^T$ 是 \mathbf{a} 和 \mathbf{b} 的负包卷积, 令 $\hat{\mathbf{a}}$, $\hat{\mathbf{b}}$ 和 $\hat{\mathbf{d}}$ 分别表示 $[a_0, \psi a_1, \dots, \psi^{n-1} a_{n-1}]^T$, $[b_0, \psi b_1, \dots, \psi^{n-1} b_{n-1}]^T$ 和 $[d_0, \psi d_1, \dots, \psi^{n-1} d_{n-1}]^T$, 那么 $\hat{\mathbf{d}} = F^{-1}(F(\hat{\mathbf{a}}) \cdot F(\hat{\mathbf{b}}))$ 。

证明: 通过观察 $\psi^n = -1$, 证明类似于定理 7.1, 详细的证明过程留作练习。 \square

7.2 快速傅里叶变换算法

如果假设环 R 中任意元素的算术操作只需要一步便可完成, 那么 R^n 中向量 \mathbf{a} 的傅里叶变换和傅里叶逆变换均只需要 $O_A(n^2)$ 时间。当 n 是 2 的幂时, 能做得更好。计算傅里叶变换和傅里叶逆变换的算法均只需要 $O_A(n \log n)$ 时间, 可推测这就是最优情况。这一节给出傅里叶变换的算法, 其逆变换算法类似, 所以留给读者自己完成。快速傅里叶变换 (FFT) 的基本思想本质是代数。我们发现 $A\mathbf{a}$ 中隐含 n 个部分和之间的相似性。本节假设 $n = 2^k$, 其中 k 是整数。

前面讲到 $A\mathbf{a}$ 的值等同于多项式 $p(x) = \sum_{j=0}^{n-1} a_j x^j$, 其中 $x = \omega^0, \omega^1, \dots, \omega^{n-1}$ 。然而, 多项式 $p(x)$ 在 $x = a$ 点上的值等于 $p(x)$ 除以 $x - a$ 的余数。[为了叙述方便, 假设 $p(x) = (x - a)q(x) + c$, 其中 c 是常量, 那么 $p(a) = c$ 。] 这样, 当一个 $(n-1)$ 次多项式 $p(x) = \sum_{j=0}^{n-1} a_j x^j$ 除以每个 $x - \omega^0, x - \omega^1, \dots, x - \omega^{n-1}$ 时, 傅里叶变换的计值归约为找余数的操作。

简单地用 $p(x)$ 除以每一个 $x - \omega^i$ 需要 $O(n^2)$ 过程。为了获得快速算法, 把 $x - \omega^i$ 两两相乘, 然后把产生的 $n/2$ 个多项式再两两相乘, 直到最后剩下两个多项式 q_1 和 q_2 , 这两个多项式的每一个都是半数的 $x - \omega^i$ 乘积。接着, 用 $p(x)$ 分别除以 q_1 和 q_2 , 得到两个余数多项式 $r_1(x)$ 和 $r_2(x)$, 它们的最大幂为 $n/2 - 1$ 。对于每个 ω^i , $x - \omega^i$ 都是 q_1 的因子, 计算 $p(x)$ 除以 $x - \omega^i$ 的余数等同于计算 $r_1(x)$ 除以 $x - \omega^i$ 的余数。类似地, 对于每个 ω^i , $x - \omega^i$ 都是 q_2 的因子。这样, 计算 $p(x)$ 除以所有 $x - \omega^i$ 的余数等于计算多项式 $r_1(x)$ 和 $r_2(x)$ 除以所有 $n/2$ 个 $x - \omega^i$ 的余数。递归使用这种分治法比直接用 $p(x)$ 逐个除以 $x - \omega^i$ 更加高效。

$x - \omega^i$ 相乘时, 期望乘积包含交叉乘项 (cross product terms)。然而, 通过合理安排 $x - \omega^i$ 的顺序, 可以使所有得到的乘积有 $x^j - \omega^i$ 这种形式, 进一步减少花费在多项式乘和除的时间。下面详细介绍这种思想。

假设 c_0, c_1, \dots, c_{n-1} 是 $\omega^0, \omega^1, \dots, \omega^{n-1}$ 的排列, 后面将说明。定义多项式 q_{lm} , 其中 $0 \leq m \leq k$, l 是 2^m 的整数倍, 且 $0 \leq l \leq 2^k - 1$, 该多项式如下:

$$a_{lm} = \prod_{j=l}^{l+2^m-1} (x - c_j)$$

这样, 多项式 q_{0k} 为积 $(x-c_0)(x-c_1)\cdots(x-c_{n-1})$, q_{0l} 是 $x-c_l$, 一般

$$q_{lm} = q_{l,m-1} q_{l+2^{m-1}, m-1}$$

存在 2^{k-m} 个第二个下标为 m 的多项式, 每一个 $x-c_l$ 正好是这些多项式中的一个因子。多项式 q_{lm} 如图 7-1 中所示 (参见 8.4 和 8.5 节)。

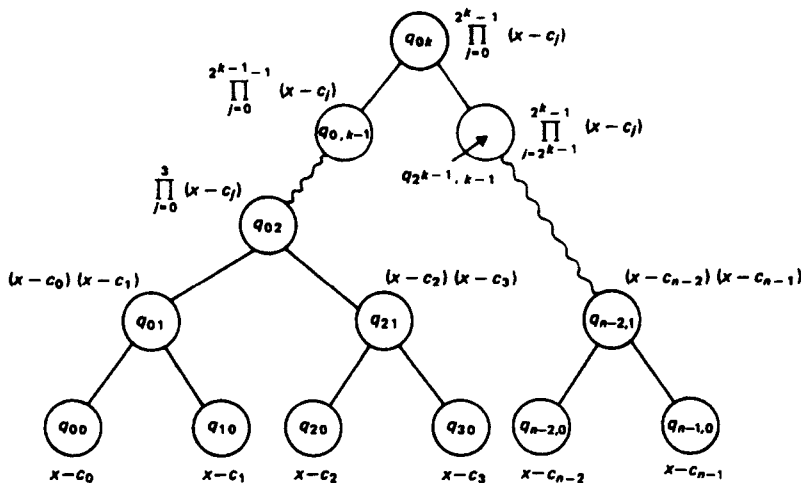


图 7-1 多项式 q_{lm}

我们是对每个 l 计算 $p(x)/q_{l0}(x)$ 的余数。为此, 首先对每个多项式 q_{lm} 计算 $p(x)/q_{lm}(x)$ 的余数, 从 $m=k-1$ 开始到 $m=0$ 结束。

假设我们已经计算了 (2^m-1) 次多项式 r_{lm} , r_{lm} 是 $p(x)/q_{lm}(x)$ 的余数。[假设 $r_{0k}=p(x)$] 因为 $q_{lm}=q'q''$, 其中 $q'=q_{l,m-1}$ 和 $q''=q_{l+2^{m-1}, m-1}$, $m-1$, 我们断言 $p(x)/q'(x)$ 和 $r_{lm}/q'(x)$ 有相同的余数, 并且这个断言对于 $q''(x)$ 也是成立的。在证明中假设

$$p(x) = h_1(x)q'(x) + r_{l,m-1}$$

其中多项式的次数 $r_{l,m-1}$ 最大可达 $2^{m-1}-1$ 。因为

$$p(x) = h_2(x)q_{lm}(x) + r_{lm}$$

我们可以得到下式

$$h_1(x)q'(x) + r_{l,m-1} = h_2(x)q_{lm}(x) + r_{lm} \quad (7-5)$$

在式(7-5)两边同时除以 $q'(x)$, 可以看到 $h_1(x)q'(x)$ 和 $h_2(x)q_{lm}(x)$ 余数为 0, 所以 $r_{lm}/q'(x)$ 的余数为 $r_{l,m-1}$ 。

这样, 通过 $(2^{k-m}-1)$ 次多项式 r_{lm} 除以 $q'(x)$ 和 $q''(x)$, 得到 $p(x)/q'(x)$ 和 $p(x)/q''(x)$ 的余数, 而不是通过 $(2^{k-m}-1)$ 次多项式 $p(x)$ 除以 $q'(x)$ 和 $q''(x)$ 得到余数。进行除法本身节省一些时间。但是, 我们可以做得更好。通过对 ω 的幂进行合理的排序 c_0, c_1, \dots, c_{n-1} , 可以保证每个多项式 q_{lm} 存在某个 s , 使之具有 $x^{2^s} - \omega^s$ 这种形式。除以这种形式的多项式是非常简单的。

引理 7.2 假设 $n=2^k$, ω 是主 n 次单位根。对于 $0 \leq j < 2^k$, 假设 $[d_0 d_1 \cdots d_{k-1}]$ 是整数 j 的二进制表示, $\text{rev}(j)$ 是二进制表示为 $[d_{k-1} d_{k-2} \cdots d_0]$ 的整数。假设 c_j 是 $\omega^{\text{rev}(j)}$, $q_{lm} = \prod_{j=l}^{l+2^m-1} (x-c_j)$,

$$\ominus j = \sum_{i=0}^{k-1} d_{k-1-i} 2^i.$$

那么 $q_{lm} = x^{2^l} - \omega^{\text{rev}(l/2^m)}$ 。

证明：通过对 m 进行归纳。归纳基础 $m=0$ ，由定义 $q_0 = x - c_l = x - \omega^{\text{rev}(l)}$ 可知，结论显然成立。由归纳步知，当 $m>0$ ，

$$\begin{aligned} q_{lm} &= q_{l,m-1} q_{l+2^{m-1}, m-1} \\ &= (x^{2^{m-1}} - \omega^{\text{rev}(l/2^{m-1})}) (x^{2^{m-1}} - \omega^{\text{rev}(l/2^{m-1}+1)}) \end{aligned}$$

其中 $l/2^{m-1}$ 是一个从 $0 \sim 2^k - 1$ 之间的偶数，那么

$$\omega^{\text{rev}(l/2^{m-1}+1)} = \omega^{2^{k-1} + \text{rev}(l/2^{m-1})} = -[\omega^{\text{rev}(l/2^{m-1})}]$$

因为 $\omega^{2^{k-1}} = \omega^{n/2} = -1$ ，那么有

$$q_{lm} = x^{2^m} - \omega^{2\text{rev}(l/2^{m-1})} = x^{2^m} - \omega^{\text{rev}(l/2^m)}$$

因为 $\text{rev}(2t) = (\frac{1}{2})\text{rev}(t)$ 。 □

例 7.1 如果 $n=8$ ，列表 c_0, c_1, \dots, c_7 ，是 $\omega^0, \omega^4, \omega^2, \omega^6, \omega^1, \omega^5, \omega^3, \omega^7$ 。 q_{lm} 如图 7-2 所示。

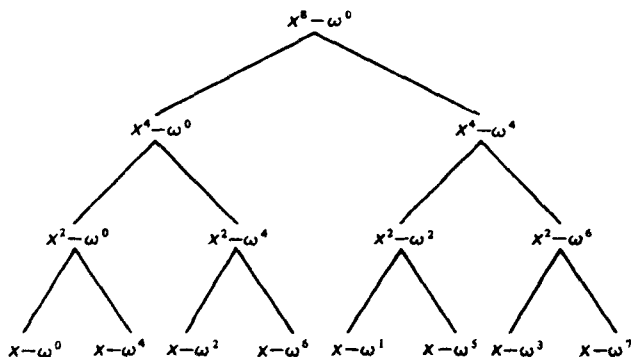


图 7-2 对引理 7.2 中 q_{lm} 的说明

q_{lm} 用来求余数如下。首先计算 $p(x)/(x^4 - \omega^0)$ 和 $p(x)/(x^4 - \omega^4)$ 的余数 r_{02} 和 r_{42} ，其中多项式 $p(x) = \sum_{j=0}^7 a_j x^j$ 。然后计算 $r_{02}/(x^2 - \omega^0)$ 和 $r_{42}/(x^2 - \omega^4)$ 的余数 r_{01} 和 r_{21} ，以及 $r_{42}/(x^2 - \omega^2)$ 和 $r_{42}/(x^2 - \omega^6)$ 的余数 r_{41} 和 r_{61} 。最后，计算 $r_{00}, r_{10}, r_{20}, \dots, r_{70}$ ，其中 r_{00} 和 r_{10} 是 $r_{01}/(x - \omega^0)$ 和 $r_{01}/(x - \omega^4)$ 的余数， r_{20} 和 r_{30} 是 $r_{21}/(x - \omega^2)$ 和 $r_{21}/(x - \omega^6)$ 的余数，如此类推。

更多这方面的例子在 8.5 节中。 □

从前面可以知道， q_{lm} 有 $x^i - c$ 这种形式现在说明 $p(x)$ 除以 $x^i - c$ 的余数是很容易计算的。

引理 7.3 假设

$$p(x) = \sum_{j=0}^{2i-1} a_j x^j$$

并且 c 是一个常数。那么 $p(x)/(x^i - c)$ 的余数为

$$r(x) = \sum_{j=0}^{i-1} (a_j + c a_{j+i}) x^j$$

证明：通过观察可以知道， $p(x)$ 可以表示为：

$$\left[\sum_{j=0}^{t-1} a_{j+t} x^j \right] (x^t - c) + r(x)$$

□

计算任意 $(2t-1)$ 次多项式除以 $x^t - c$ 的余数都可以在 $O_A(t)$ 步内完成，该方法比任意已知计算 $(2t-1)$ 次多项式除以任意 t 次多项式的余数都快。

完整的FFT算法如下。

算法 7.1 快速傅里叶变换。

输入：向量 $\mathbf{a} = [a_0, a_1, \dots, a_{n-1}]^T$ ，其中 $n=2^k$ ，且 k 为整数。

输出：向量 $F(\mathbf{a}) = [b_0, b_1, \dots, b_{n-1}]^T$ ，其中 $b_i = \sum_{j=0}^{n-1} a_j \omega^{ij} (0 \leq i < n)$ 。

方法：如图 7-3 所示。

□

```

begin
1.  let  $r_{0k}$  be  $\sum_{j=0}^{n-1} a_j x^j$ ;
    comment 一个多项式用系数表示, 所以第1行没有涉及计算. 当  $\sum_{j=0}^{n-1} a_j x^j$  除以
            $q_{lm}$  时,  $r_{lm}$  表示余数;
2.  for  $m \leftarrow k-1$  step  $-1$  until  $0$  do
3.    for  $l \leftarrow 0$  step  $2^{m+1}$  until  $n-1$  do
        begin
4.          let  $r_{l,m+1}$  be  $\sum_{j=0}^{2^m-1} a_{l+j} x^j$ ;
            comment 系数  $r_{l,m+1}$  计算低次的余数
5.           $s \leftarrow \text{rev}(l/2^m)$ ;
6.           $r_{lm}(x) \leftarrow \sum_{j=0}^{2^m-1} (a_j + \omega^s a_{j+x^m}) x^j$ ;
7.           $r_{l+x^m,m}(x) \leftarrow \sum_{j=0}^{2^m-1} (a_j + \omega^{s+n/2} a_{j+x^m}) x^j$ 
            end;
8.    for  $l \leftarrow 0$  until  $n-1$  do  $b_{\text{rev}(l)} \leftarrow r_{lm}$ 
end

```

图 7-3 快速傅里叶变换

为了修改算法 7.1 来计算逆变换，仅通过修改第 6、7 行中 ω 的指数符号，把 ω 替换为 ω^{-1} ，然后在第 8 行中用 n 除 $b_{\text{rev}(l)}$ 。

例 7.2 假设 $n=8$ ，那么 $k=3$ 。对于 $m=2$ ，行 3~7 的循环操作只有当 $l=0$ 时执行。在第 4 行， r_{03} 为 $\sum_{j=0}^7 a_j x^j$ ，在第 5 行， s 为 0，在第 6、7 行有

$$r_{02} = (a_3 + a_7)x^3 + (a_2 + a_6)x^2 + (a_1 + a_5)x + (a_0 + a_4)$$

和

$$r_{42} = (a_3 + \omega^4 a_7)x^3 + (a_2 + \omega^4 a_6)x^2 + (a_1 + \omega^4 a_5)x + (a_0 + \omega^4 a_4)$$

对于 $m=1$ ， l 在 $0 \sim 4$ 之间取值。当 $l=0$ 时，在第 5 行有 $s=0$ 。这样在执行 6、7 行后 $r_{01} = (a_1 + a_3 + a_5 + a_7)x + (a_0 + a_2 + a_4 + a_6)$

和

$$r_{21} = (a_1 + \omega^4 a_3 + a_5 + \omega^4 a_7)x + (a_0 + \omega^4 a_2 + a_4 + \omega^4 a_6)$$

当 $l=4$ 时，计算 $s=2$ 。可以通过第 6、7 两行和上面的公式求出的 r_{42} 得到

$$r_{41} = (a_1 + \omega^2 a_3 + \omega^4 a_5 + \omega^6 a_7)x + (a_0 + \omega^2 a_2 + \omega^4 a_4 + \omega^6 a_6)$$

$$r_{61} = (a_1 + \omega^6 a_3 + \omega^4 a_5 + \omega^2 a_7)x + (a_0 + \omega^6 a_2 + \omega^4 a_4 + \omega^2 a_6)$$

最后, 对于 $m=0$, l 取值 0、2、4 和 6。例如, 对于 $l=4$, 有 $s=1$, 由 r_{41} 计算得到:

$$r_{40} = a_0 + \omega a_1 + \omega^2 a_2 + \cdots + \omega^7 a_7$$

在执行到第 8 行的 for 循环前, r_{40} 总是次为 0 的多项式, 也就是常数。例如, 当 $l=4$, $\text{rev}(l)=1$, r_{40} 赋值为 b_1 。对于 b_1 的这个公式总是符合 b_l 的定义。□

下面证明算法 7.1 的正确性。

定理 7.3 算法 7.1 计算离散傅里叶变换。

证明: 在第 6 行, $r_{lm} = r_{l,m+1}/q_{lm}$, 第 7 行, $r_{l+2,m} = r_{l,m+1}/q_{l+2,m}$ 。通过使用引理 7.2 和 7.3, 很容易通过对 $k-m$ 归纳证明 $r_{lm} = \sum_{j=0}^{n-1} a_j x^j$ 除以 q_{lm} 的余数。那么, 对于 $m=0$, 引理 7.3 保证第 8 行的 for 循环给每个 b_l 赋予正确的(常量)余数。□

定理 7.4 算法 7.1 需要的时间为 $O_A(n \log n)$ 。

证明: 每次执行第 6 行和第 7 行需要 $O_A(2^m)$ 步。对于固定的 m , 第 3~7 行的循环迭代 $n/2^{m+1}$ 次, 需要的总时间与 m 无关, 是 $O_A(n)$ 。从第 2 行开始的外层循环总共执行 $\log n$ 次, 所以总代价是 $O_A(n \log n)$ 。第 8 行的循环实际上不需要算数运算。□

我们已经在定理 7.4 中假设 n 是固定的。这样, 在 5~8 行, ω 的幂、 s 和 $\text{rev}(l)$ 值都可以由 l 先计算, 在执行时作为常量使用。如果程序中把 n 作为一个参数, 那么仍旧可以利用 $O(n)$ 步计算出 ω 的幂, 并把它存储在 RAM 表中。而且, 即使需要在第 5~8 行花费 $O(\log n)$ 步由 l 计算 s 和 $\text{rev}(l)$, 计算不会超过 $3n$ 步。所以在 RAM 中整个算法的时间复杂度为 $O(n \log n)$ 。

推论 1 可以在 $O_A(n \log n)$ 步内计算出 $\mathbf{a} \otimes \mathbf{b}$, 其中 \mathbf{a} 、 \mathbf{b} 是长度为 n 的向量。

证明: 通过定理 7.1、7.3 和 7.4 直接可以得到。□

推论 2 可以在 $O_A(n \log n)$ 步内计算出 \mathbf{a} 、 \mathbf{b} 向量的正负包卷积。

FFT 的算法 7.1 提供了算法开发的直观叙述。如果要实际计算傅里叶变换, 那么只需对其系数进行处理, 对算法进行一定的简化。具体的工作如算法 7.2 中所示。

算法 7.2 简化的傅里叶变换算法(simplified FFT algorithm)。

输入: 向量 $\mathbf{a} = [a_0, a_1, \dots, a_{n-1}]^T$, 其中 $n=2^k$ 且 k 为整数。

输出: 向量 $F(\mathbf{a}) = [b_0, b_1, \dots, b_{n-1}]^T$, 其中 $b_i = \sum_{j=0}^{n-1} a_j \omega^{ij} (0 \leq i < n)$ 。

方法: 利用图 7-4 中的程序。为了概念上简化, 在这个程序中用临时数组 S 存储前一步的值。实际上, 计算在原地完成。□

第 3 行第一次执行时, 多项式 $p(x) = \sum_{i=1}^{n-1} a_i x^i$ 的系数存储在数组 S 中。第 6 行第一次执行时, $p(x)$ 除以 $x^{n/2} - 1$ 和 $x^{n/2} - \omega^{n/2}$, 余数分别为:

$$\sum_{i=0}^{n/2-1} (a_i + a_{i+n/2}) x^i \text{ 和 } \sum_{i=0}^{n/2-1} (a_i + \omega^{n/2} a_{i+n/2}) x^i$$

在第 2 次执行第 3 行时, 两个余数的系数存储在数组 S 中。第 1 个余数的系数存储在数组 S 的前半部分, 第 2 个余数的系数则存储在数组 S 的后半部分。

```

begin
1.   for i ← 0 until 2k - 1 do R[i] ← ai;
2.   for l ← 0 until k - 1 do
       begin
3.       for i ← 0 until 2k - 1 do S[i] ← R[i];
4.       for i ← 0 until 2k - 1 do
           begin
5.           设 [d0d1...dk-1] 是整数 i 的二进制表示;
6.           R[[d0...dk-1]] ←
               S[[d0...dl-10dl+1...dk-1]]
               + ω[dldl-1...d00]S[[d0...dl-11dl+1...dk-1]]
           end
       end
7.   for i ← 0 until 2k - 1 do
       b[d0...dk-1]] ← R[[dk-1...d0]]
   end
end

```

图 7-4 简化的傅里叶变换

对第 6 行的第二次执行, 两个余数多项式除以两个形为 $x^{n/4} - \omega'$ 的多项式。这产生了 4 个次数为 $n/4 - 1$ 的余数。第 3 行的第三次执行把这四个余数的系数存储到 S 中, 其余类似。第 7 行重新对结果的分量按照正确的顺序进行重新排列。在计算 $x - \omega'$ 的乘积时, 为了消去叉乘项而要求置换单位根, 所以这一行操作是必需的。 $n = 8$ 的部分过程在图 7-5 中进行了说明。

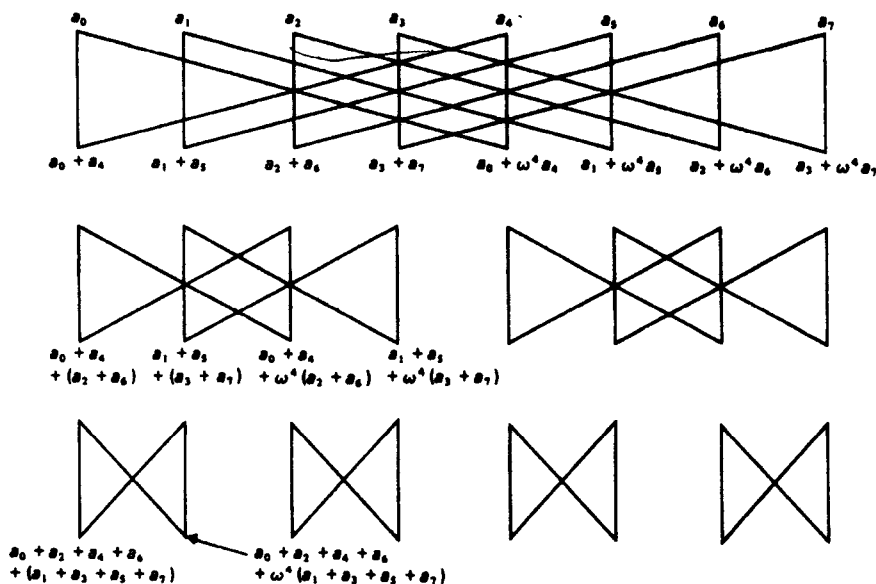


图 7-5 通过算法 7.2 计算 FFT 的过程解释。因为缺少空间, 省略了一些余数多项式的系数

7.3 使用位操作的 FFT

在使用傅里叶变换简化卷积计算的应用中, 时常需要确切的结果。如果计算是在实数环上, 那么必须把实数转换成精度有限的数, 导致了错误的产生。这些错误可以通过在一个有限域中

进行计算而避免[⊖]。例如,为了计算 $\mathbf{a} = [a_0, a_1, a_2, 0, 0]^T$ 和 $\mathbf{b} = [b_0, b_1, b_2, 0, 0]^T$ 的卷积,可以令2作为第5次单位根应用于计算模31。然后变换 \mathbf{a} 和 \mathbf{b} , 执行两两相乘。并且在 $\mathbf{a} \otimes \mathbf{b} \bmod 31$ 计算逆变换的精确结果[⊕]。使用有限域的困难在于找到一个合适的第 n 次单位根的域。作为替换方法,可以使用模整数 m 的环 R_m , 选择 m 值使环存在第 n 次单位根 ω [⊗]。对于给定的 n , 并不能马上找到 ω 和 m , 使 ω 是模 m 的整数环的第 n 次单位根。而且,也不能把 m 定义为太大的值,由此造成模 m 计算代价的增加,而不被允许。幸运的是,如果 n 是2的幂,那么存在一个合适的 m , 近似为 2^n 。特别是定理7.5分析了当 n 和 ω ($\omega > 1$) 都是2的幂时,可以通过傅里叶变换、分量乘法(componentwise multiplication)和逆变换在模为 $(\omega^{n/2} + 1)$ 的整数环中计算卷积。首先,引理7.4和7.5将给出预备结果。在这些引理中,假设 $R = (S, +, \cdot, 0, 1)$ 是一个交换环,其中 $n = 2^k$, $k \geq 1$ 。

引理7.4 对于所有的 $a \in S$,

$$\sum_{i=0}^{n-1} a^i = \prod_{i=0}^{k-1} (1 + a^{2^i})$$

证明:对 k 进行归纳。归纳基础 $k=1$ 时,结果显然。现在,观察下式

$$\sum_{i=0}^{n-1} a^i = (1+a) \sum_{i=0}^{n/2-1} (a^2)^i \quad (7-6)$$

通过规约假设和把 a 代换为 a^2 , 可以得到

$$\sum_{i=0}^{n/2-1} (a^2)^i = \prod_{i=1}^{k-2} [1 + (a^2)^{2^i}] = \prod_{i=1}^{k-1} [1 + a^{2^i}] \quad (7-7)$$

利用式(7-7)代替式(7-6)的右边,归纳假设成立。□

引理7.5 假设 $m = \omega^{n/2} + 1$, 其中 $\omega \in S$, $\omega \neq 0$ 。那么对于 $1 \leq p < n$, 有 $\sum_{i=0}^{n-1} \omega^{ip} \equiv 0 \bmod m$ 。

证明:通过引理7.4可以证明存在 j ($0 \leq j < k$), $1 + \omega^{2^j} \equiv 0 \bmod m$ 。假设 $p = 2^j p'$, 其中 p' 是奇数, 那么 $0 \leq s < k$ 。选择 j , 使得 $j+s = k-1$ 。那么 $1 + \omega^{2^j} = 1 + \omega^{2^{j+s} p'} = 1 + (m-1)^{p'}$ 。但是 $(m-1) \equiv -1 \bmod m$, 并且 p' 是奇数, 所以 $(m-1)^{p'} \equiv -1 \bmod m$ 。对于 $j = k-1-s$, $1 + \omega^{2^j} \equiv 0 \bmod m$ 成立。□

定理7.5 假设 n 和 ω 是2的幂, $m = \omega^{n/2} + 1$ 。假设 R_m 是模 m 的整数环, 那么在 R_m 中, n 有一个乘法逆模 m (multiplicative inverse modulo m), 并且 ω 是一个主 n 次单位根。

证明:因为 n 是2的幂, m 是奇数, 所以有 m 和 n 互质。这样, n 有一个乘法逆模 m [⊗]。因为 $\omega \neq 1$, $\omega^n = \omega^{n/2} \omega^{n/2} \equiv (-1)(-1) = 1 \bmod (\omega^{n/2} + 1)$ 。然后从引理7.5中可以得出 ω 是 R_m 中的一个主 n 次单位根。□

定理7.5的重要性在于卷积定理在模 $2^{n/2} + 1$ 的整数环上是有效的。如果希望计算两个 n 维整数向量的卷积, 并且卷积的分量在0到 $2^{n/2}$ 之间, 那么肯定可以得到确切的答案。如果卷积的分量不在0到 $2^{n/2}$ 之间, 那么模 $2^{n/2} + 1$ 是正确的。

现在几乎可以得出计算卷积时模 m 所需要的位操作的数目。首先考虑在计算整数模 m 的余数中用到的位操作的数目, 这是从模 m 的算术模数推出位操作数必不可少的。

假设存在整数 p 使 $m = \omega^p + 1$ 。使用“除9余数核对”(以9除得到的余数, 再以相同的除法,

⊖ “域”的定义见12.1节。

⊗ 当然必须确定结果元素在0~30之间, 否则不能恢复它们。通常, 必须选择足够大的模来使其能恢复。

⊕ 目前为止, 如果把 ω 当作复数 $e^{2\pi i/n}$, 并且像在复数域上计算一样, 那么没有误入歧途。但 ω 必须当作整数, 所有计算必须在模 m 的有限整数环上。

⊗ 该结论是数论的基本定理。在8.8节中, 可以看到, 如果 a 和 b 互质, 那么存在整数 x 和 y , 使 $ax + by = 1$ 成立。即 $ax \equiv 1 \bmod b$ 。假设 $b = m$, a 是 n 可以得到结论。

以9除之。——译者注)技术的推广来计算 a 模 m 的值。如果基于 ω^p 记法, a 可以表示为 p 位的 l 块序列,那么 a 模 m 可以通过交替加减 p 位 l 块的方法进行计算。

引理 7.6 假设 $m = \omega^p + 1$, 且 $a = \sum_{i=0}^{l-1} a_i \omega^{pi}$, 其中对于每个 i , 有 $0 \leq a_i < \omega^p$ 。那么, $a \equiv \sum_{i=0}^{l-1} a_i (-1)^i$ 模 m 。

证明: 观察 $\omega^p \equiv -1$ 模 m 。 □

注意: 引理 7.6 中的(块数) l 固定不变, 那么计算 a 模 m 的余数可以通过 $O_b(p \log \omega)$ 次位操作实现。

例 7.3 假设 $n=4$, $\omega=2$, $m=2^2+1$ 。那么在引理 7.6 中 $p=2$ 。考虑二进制数 $a=101100$ 。这里 $a_0=00$, $a_1=11$, $a_2=10$ 。我们可以计算 $a_0 - a_1 + a_2 = -1$, 然后加 m , 可以发现 $a \equiv 4$ 模 5。因为 a 为 44, 所以结果得到验证。 □

引理 7.6 提供了一个计算 a 模 m 的有效方法。它在下面的定理中起到了非常重要的作用, 下面这个定理给出了在计算离散傅里叶变换及其逆变换计算时所需要的位操作数的上界。

定理 7.6 假设 ω 和 n 是 2 的幂, $m = \omega^{n/2} + 1$ 。假设 $[a_0, a_1, \dots, a_{n-1}]^T$ 是整数向量, 对于每个 i , 有 $0 \leq a_i < m$ 。那么 $[a_0, a_1, \dots, a_{n-1}]^T$ 的离散傅里叶变换及其逆可以在 $O_b(n^2 \log n \log \omega)$ 步内计算模 m 。

证明: 利用算法 7.1 或者 7.2。在逆变换中, 用 ω^{-1} 置换 ω 并且用 n^{-1} 乘以每个结果。整数模 m 可以被一个长度为 $b = ((n/2) \log \omega) + 1$ 位的字符串表示。由于 $m = 2^{b-1} + 1$, 模 m 的余数可以通过位串 $00 \dots 0 \sim 100 \dots 0$ 表示。

算法 7.1 要求与 ω 的幂进行模 m 的整数加法和模 m 的乘法操作。这些操作执行 $O(n \log n)$ 次。利用引理 7.6, 模 m 的加法要求 $O_b(b)$ 步, 其中 $b = ((n/2) \log \omega) + 1$ 。因为 ω 是 2 的幂, 所以乘以 ω^p 等于左移 $p \log \omega$ 位, 其中 $0 \leq p < n$ 。得到的整数至多有 $3b-2$ 位, 通过引理 7.6, 位移和计算余数要求 $O_b(b)$ 步。如此, 前向傅里叶变换的时间复杂度为 $O_b(bn \log n)$, 也就是 $O_b(n^2 \log n \log \omega)$ 。

逆变换要求乘以 ω^{-p} 和 n^{-1} 。因为 $\omega^p \omega^{n-p} \equiv 1$ 模 m , 所以有 $\omega^{n-p} \equiv \omega^{-p}$ 模 m 。那么乘以 ω^{-p} 的结果和乘以 ω^{n-p} 得到的结果一致。后者是左位移 $(n-p) \log \omega$ 位, 得到的整数至多 $3b-2$ 位。同样, 利用引理 7.6 可以在 $O_b(b)$ 步内得到余数。最后, 考虑与 n^{-1} 相乘。如果 $n=2^k$, 那么左位移 $n \log \omega - k$ 位, 保证该数至多 $3b-2$ 位, 然后, 通过引理 7.6 计算余数。这样, 逆变换也要求 $O_b(n^2 \log n \log \omega)$ 步。 □

例 7.4 假设 $\omega=2$, $n=4$, $m=5$ 。对向量 $[a_0, a_1, a_2, a_3]^T$ 进行傅里叶变换, 其中 $a_i = i$ 。由于对所有的 i , 有 $a_i < 5$, 如果计算模 m , 期望恢复这个向量。用三位来表示数, 除了临时结果, 只有 000, \dots , 100 会实际用到。

在算法 7.1 中, 必须计算图 7-6 所示的多项式系数, 其中 2 已经被 ω 替换了。

a_0	a_1	a_2	a_3
$a_0 + a_2$	$a_1 + a_3$	$a_0 + 4a_2$	$a_1 + 4a_3$
$a_0 + a_2 + (a_1 + a_3)$ $= a_0 + a_1 + a_2 + a_3$ $= b_0$	$a_0 + a_2 + 4(a_1 + a_3)$ $= a_0 + 4a_1 + a_2 + 4a_3$ $= b_2$	$a_0 + 4a_2 + 2(a_1 + 4a_3)$ $= a_0 + 2a_1 + 4a_2 + 8a_3$ $= b_1$	$a_0 + 4a_2 + 8(a_1 + 4a_3)$ $= a_0 + 8a_1 + 4a_2 + 2a_3$ $= b_3$

图 7-6 当 $n=4$ 时, 计算快速傅里叶变换

图 7-6 中的实际值为:

$$\begin{array}{cccc}
 a_0 = 000 & a_1 = 001 & a_2 = 010 & a_3 = 011 \\
 a_0 + a_2 = 010 & a_1 + a_3 = 100 & a_0 + 4a_2 = 011 & a_1 + 4a_3 = 011 \\
 b_0 = 001 & b_1 = 011 & b_2 = 100 & b_3 = 010
 \end{array}$$

变换 $[0, 1, 2, 3]^T$ 是 $[1, 4, 3, 2]^T$ 模 5。考虑最后一行的最后一个元素 b_3 。它由中间行的最后两个元素通过下面的步骤计算得到。

$$\begin{array}{ll}
 \text{计算 } a_1 + 4a_3 & 011 \\
 \text{向左移三位(乘以 8)} & 11000 \\
 \text{分为长度为 2 的三块} & 1 \ 10 \ 00 \\
 \text{把第一块和第三块相加, 再减去第二块} & -1 \\
 \text{加上 } m=5 & 100 \\
 \text{加上 } a_0 + 4a_2 = 011 & 111 \\
 \text{减去 } m & 010
 \end{array}$$

通过反转操作, 可以得到 $2^{-1} \equiv 8$ 模 5, $4^{-1} \equiv 4$ 模 5 和 $8^{-1} \equiv 2$ 模 5。如此, 逆变换公式可以从图 7-6 中通过交换 a_i 和 b_i 以及 2 和 8 得到。计算过程如下:

$$\begin{array}{cccc}
 b_0 = 001 & b_1 = 100 & b_2 = 011 & b_3 = 010 \\
 b_0 + b_2 = 100 & b_1 + b_3 = 001 & b_0 + 4b_2 = 011 & b_1 + 4b_3 = 010 \\
 4a_0 = 000 & 4a_2 = 011 & 4a_1 = 100 & 4a_3 = 010
 \end{array}$$

最后, 把每个结果除以 4 (乘以 4, 因为有 $4^{-1} \equiv 4$ 模 5), 可以由 $[a_0, a_1, a_2, a_3]^T$ 得到 $[0, 1, 2, 3]^T$ 。□

定理 7.6 的推论 假设 $O_b(M(k))$ 是计算两个 k 位整数的积所需要的步数。假设 a 和 b 是长度为 n 的向量, 且该向量的整数分量在 $0 \sim \omega^n$ 范围内, n 和 ω 是 2 的幂。可以在 $O_b(\text{MAX}[n^2 \log n \log \omega, nM(n \log \omega)])$ 时间内计算 $a \otimes b$, 或者 a 和 b 模 $\omega^n + 1$ 的正包卷积或负包卷积。

定理 7.6 的推论中函数的第一项表示变换的时间, 第二项表示做 $2n$ 次 $(n \log \omega + 1)$ 位整数乘法所需要的代价。最好的情况下, $M(k)$ 的值可以是 $k \log k \log \log k$ (见 7.5 节)。在这种情况下, 第二项支配第一项, 所以该卷积操作需要 $O_b(n^2 \log n \log \log n \log \omega \log \log \omega \log \log \omega)$ 步。

7.4 多项式乘积

计算两个一元多项式的乘积问题实际上和计算两个序列数的卷积操作相同。也就是说

$$\left(\sum_{i=0}^{n-1} a_i x^i \right) \left(\sum_{j=0}^{n-1} b_j x^j \right) = \sum_{k=0}^{2n-2} c_k x^k, \quad \text{其中 } c_k = \sum_{m=0}^{n-1} a_m b_{k-m}$$

和前面的一样, 如果 $p < 0$ 或者 $p \geq n$, 那么 a_p 和 b_p 取值为 0。并且 c_{2n-1} 必须为 0。下面是定理 7.4 的一些推论。

定理 7.4 的推论 3 两个 n 次多项式乘积的系数可以在 $O_\lambda(n \log n)$ 步内计算得到。

证明: 从定理 7.4 的推论 1 和上面的观察可以直接得到。□

定理 7.4 的推论 4 假设可以在 $M(k)$ 步内计算出两个 k 位整数的乘积,

$$p(x) = \sum_{i=0}^{n-1} a_i x^i \text{ 和 } q(x) = \sum_{j=0}^{n-1} b_j x^j$$

且对于所有的 i 和 j , a_i 和 b_j 是在 $0 \sim \omega^{n/2}/\sqrt{n}$ 之间的整数, 其中 n 和 ω 是 2 的幂。那么, 可以在 $O_B(\text{MAX}[n^2 \log n \log \omega, nM(n \log \omega)])$ 步内计算 $p(x)q(x)$ 的系数。

证明: 从定理 7.4 和定理 7.6 的推论可得。□

注意, 在推论 4 中, 第二项占支配位置。

实际上, 定理 7.1 有下面的解释。假设 $p(x)$ 和 $q(x)$ 是 $(n-1)$ 次多项式。在 $2n-1$ 或者更多点对 p 和 q 计值, 也就是 c_0, c_1, \dots , 然后计算 $p(c_j)q(c_j)$ 的值获得 pq 在这些点上的值。唯一的 $(2n-2)$ 次多项式可以通过这些点插值。这个 $(2n-2)$ 次多项式就是 $p(x)$ 和 $q(x)$ 的乘积。

当把傅里叶变换运用于卷积(或者等价的多项式乘法)时, 选择 $c_j = \omega^j$, 其中 ω 是主 $2n$ 次单位根。计算 p 和 q 的傅里叶变换, 也就是对它们在 c_0, c_1, \dots 点上计值。然后计算变换的两两对应乘积, 也就是 $p(c_j)$ 和 $q(c_j)$ 相乘得到 c_j 的值。接着, 运用逆变换计算 pq 。引理 7.1 保证了逆变换实际上是一个插值公式。也就是说, 实际上是通过在 $\omega^0, \omega^1, \dots, \omega^{2n-1}$ 这些点上的值恢复多项式。

7.5 Schönhage-Strassen 整数相乘算法

现在转到卷积定理的一个重要应用, 快速按位整数乘法算法(a fast bitwise integer - multiplication algorithm)。在 2.6 节中介绍到了如何通过把两个 n 位二进制整数中的每一个分为两个 $(n/2)$ 位整数在 $O(n^{\log_2 3})$ 步内计算它们的乘积。这个方法可以推广为把每个整数分为 b 块 l 位的数。通过把 b 块当做多项式的系数可以对式(2-4)中的表达式进行类似的改进。为了找到多项式乘积的系数, 通过对某个合适的点集、乘以样本值和插值来计算多项式。通过选择主单位根为计值点, 可以利用傅里叶变换和卷积定理。通过假设 b 为一个 n 的递归函数, 可以在 $O_B(n \log n \log \log n)$ 步内计算出两个 n 位数的乘法。

为了简化讨论, 把 n 限制为 2 的幂。对于 n 不是 2 的幂的情况, 可以通过加上合适数量的前导 0, 使得 n 为 2 的幂(这仅增加了常数因子)。更进一步, 我们将计算两个 n 位整数模 $(2^n + 1)$ 的乘积。为了得到两个 n 位整数准确的乘积, 必须引入前导 0 和 $2n$ 位整数模 $(2^{2n} + 1)$ 乘法, 从而再次由于常量因子增加了需要的时间。

假设 u 和 v 是要进行模 $(2^n + 1)$ 乘法的 $0 \sim 2n$ 之间的二进制整数。通过观察可知, 在二进制表示中, 2^n 需要 $n+1$ 位。如果 u 或者 v 等于 2^n , 那么用一个特殊符号 -1 来表示, 这种情况下的乘法操作很容易作为一个特例处理。如果 $u = 2^n$, 那么 uv 模 $(2^n + 1)$ 可以通过计算 $(2^n + 1 - v)$ 模 $(2^n + 1)$ 得到。

假设 $n = 2^k$, 且当 k 为偶数时 $b = 2^{k/2}$, 其他情况下 $b = 2^{(k-1)/2}$ 。假设 $l = n/b$ 。观察 $l \geq b$ 和 b 除以 l 。第一步是把 u 和 v 分为 b 块, 每块大小为 l 位。如此

$$u = u_{b-1}2^{(b-1)l} + \dots + u_12^l + u_0$$

和

$$v = v_{b-1}2^{(b-1)l} + \dots + v_12^l + v_0$$

u 和 v 的乘积由下式得到

$$uv = y_{2b-2}2^{(2b-2)l} + \dots + y_12^l + y_0 \quad (7-8)$$

其中

$$y_i = \sum_{j=0}^{b-1} u_j v_{i-j}, 0 \leq i < 2b$$

(对于 $j < 0$ 或者 $j > b-1$, 取 $u_j = v_j = 0$ 。项 y_{2b-1} 为 0, 仅是为了满足对称。)

乘积 uv 可以通过卷积定理计算。利用傅里叶变换乘法(multiplying the fourier transforms)需要 $2b$ 次乘法操作。使用卷积可以把乘法操作的次数减少到 b 。原因在于我们是在模 $2^n + 1$ 下计算 uv

的。因为 $bl = n$, 所以有 $2^n \equiv -1 \pmod{2^n + 1}$ 。这样, 通过式(7-8)和引理 7-6, 计算 $wv \pmod{2^n + 1}$ 为

$$wv \equiv (w_{b-1}2^{(b-1)l} + \cdots + w_l2^l + w_0) \pmod{2^n + 1}$$

其中, $w_i = y_i - y_{b+i}$, $0 \leq i < b$ 。

由于两个 l 位数的乘积必定少于 2^{2l} , 由于 y_i 和 y_{b+i} 分别是 $i+1$ 和 $b-(i+1)$ 位这类乘积的和, $w_i = y_i - y_{b+i}$ 必须在 $-(b-1-i)2^{2l} < w_i < (i+1)2^{2l}$ 范围内。如此, w_i 总共有最多 $b2^{2l}$ 种可能的取值。如果能够计算 $w_i \pmod{b2^{2l}}$, 那么通过加上 bw_i 和适当的位移就可以在 $O(b \log(b2^{2l}))$ 额外步内计算出 $wv \pmod{2^n + 1}$ 。

为了计算 $w_i \pmod{b2^{2l}}$, 计算 w_i 两次。一次是模 b , 还有一次是模 $2^{2l} + 1$ 。假设 w'_i 是 $w_i \pmod{b}$, \hat{w}_i 是 $w_i \pmod{2^{2l} + 1}$ 。因为 b 是 2 的幂, $2^{2l} + 1$ 是奇数, b 和 $2^{2l} + 1$ 互质。所以 w_i 可以由 w'_i 和 \hat{w}_i 通过如下公式计算。

$$w_i = (2^{2l} + 1)((w'_i - \hat{w}_i) \pmod{b}) + \hat{w}_i \pmod{b2^{2l}}$$

w_i 在 $(b-1-i)2^{2l}$ 和 $(i+1)2^{2l}$ 之间。对于每个 w_i , 从 w'_i 和 \hat{w}_i 计算 w_i 需要 $O(l + \log b)$, 总的需要时间为 $O(bl + b \log b)$ 或者 $O(n)$ 。

通过设置 $u'_i = u_i \pmod{b}$, $v'_i = v_i \pmod{b}$, 计算 $w_i \pmod{b}$ 的值, 并且生成两个 $(3b \log b)$ 位的数 \hat{u} 和 \hat{v} 如图 7-7。利用第 2.6 节中的算法计算积 $\hat{u}\hat{v}$ 需要至多 $O((3b \log b)^{1.6})$ 步, 也就是少于 $O(n)$ 步。可以看到 $\hat{u}\hat{v} = \sum_{i=0}^{2b-1} y'_i 2^{(3 \log b)i}$, 其中 $y'_i = \sum_{j=0}^{2b-1} u'_j v'_{i-j}$ 。由于 $y'_i < 2^{3 \log b}$, 所以 y'_i 可以很容易地从 $\hat{u}\hat{v}$ 中复原。 $w_i \pmod{b}$ 的值可以通过计算 $(y'_i - y'_{b+i}) \pmod{b}$ 找到。

$$\begin{array}{l} \hat{u} = u'_{b-1} 00 \dots 0 u'_{b-2} 00 \dots 0 u'_{b-3} \dots 00 \dots 0 u'_0 \\ \hat{v} = v'_{b-1} 00 \dots 0 v'_{b-2} 00 \dots 0 v'_{b-3} \dots 00 \dots 0 v'_0 \end{array}$$

图 7-7 在计算 $w_i \pmod{b}$ 中用到的复数。每一 0 块中都有 $2 \log b$ 个 0

计算完 $w_i \pmod{b}$ 后, 通过包卷积计算 $w_i \pmod{2^{2l} + 1}$ 。这涉及傅里叶变换、两两相乘和逆变换。假设 $\omega = 2^{4l/b}$, $m = 2^{2l} + 1$ 。由定理 7.5, ω 和 b 有模 m 的乘法逆 (multiplicative inverses modulo m), ω 是主 b 次单位根。如此, 在 $\psi = 2^{2l/b}$ 时 (ψ 是 $2b$ 次单位根), $[u_0, \psi u_1, \dots, \psi^{b-1} u_{b-1}]$ 和 $[v_0, \psi v_1, \dots, \psi^{b-1} v_{b-1}]$ 的负包卷积是:

$$[(y_0 - y_b), \psi(y_1 - y_{b+1}), \dots, \psi^{b-1}(y_{b-1} - y_{2b-1})] \pmod{2^{2l} + 1},$$

其中 $y_i = \sum_{j=0}^{b-1} u_j u_{i-j}$ ($0 \leq i \leq 2b-1$)。 $w_i \pmod{2^{2l} + 1}$ 可以通过恰当的位移操作得到。完整的算法概括如下。

算法 7.3 Schönhage-Strassen 整数乘法算法。

输入: 两个 n 位整数 u 和 v , 其中 $n = 2^k$ 。

输出: $(n+1)$ 位的 u 和 v 模 $2^n + 1$ 乘积。

方法: 如果 n 很小, 那么可以挑选所喜欢的算法来计算 u 和 v 模 $(2^n + 1)$ 的乘积。当 n 较大时, 如果 k 是偶数则令 $b = 2^{k/2}$, 其他情况则令 $b = 2^{(k-1)/2}$ 。假设 $l = n/b$ 。表达式 $u = \sum_{i=0}^{b-1} u_i 2^{il}$, $v = \sum_{i=0}^{b-1} v_i 2^{il}$, 其中 u_i 和 v_i 是介于 $0 \sim 2^l - 1$ 之间的整数 (也就是说 u_i 是 u 的 l 位长度的块, v_i 是 v 的 l 位长度的块)。

⊙ 如果 p_1 和 p_2 是互质的, $w \equiv q_1 \pmod{p_1}$, $w \equiv q_2 \pmod{p_2}$, 且 $0 \leq w < p_1 p_2$, 则 $w = p_2(p_2^{-1} \pmod{p_1})(q_1 - q_2 \pmod{p_1}) + q_2$ 。假设 $p_1 = b$ 并且 $p_2 = 2^{2l} + 1$ 。由于 b 是 2 的幂且 $b \leq 2^{2l}$, b 除以 2^{2l} , 那么 $2^{2l} + 1$ 模 b 的乘法逆为 1。

1. 计算模 $2^{2l} + 1$ 的 $[u_0, \psi u_1, \dots, \psi^{b-1} u_{b-1}]^T$ 和 $[v_0, \psi v_1, \dots, \psi^{b-1} v_{b-1}]^T$ 的傅里叶变换, 其中 $\psi = 2^{2l/b}$, ψ^2 是主 b 次单位根。

2. 使用递归计算每两两相乘的算法 7.3 计算步 1 中模 $2^{2l} + 1$ 的傅里叶变换的两两乘积。(在这里把 2^{2l} 当作特例来容易地处理。)

3. 对第 2 步中两两相乘的向量计算模 $2^{2l} + 1$ 的傅里叶逆变换。这个计算结果是 $[w_0, \psi w_1, \dots, \psi^{b-1} w_{b-1}]^T$ 模 $2^{2l} + 1$, 其中 $w_i = [u_0, u_1, \dots, u_{b-1}]^T$ 和 $[v_0, v_1, \dots, v_{b-1}]^T$ 负包卷积的第 i 项。通过 $\psi^i w_i$ 和 ψ^{-i} 相乘模 $2^{2l} + 1$ 可以计算出 $w'_i = w_i$ 模 $2^{2l} + 1$ 。

4. 按如下方式计算 $w'_i = w_i \bmod b$:

a) 假设 $u'_i = u_i$ 模 b , $v'_i = v_i$ 模 b , 其中 $0 \leq i < b$;

b) 通过串 u'_i 、 v'_i 和插入 $2 \log b$ 个 0 来构造数 \hat{u} 和 \hat{v} 。也就是 $\hat{u} = \sum_{i=0}^{b-1} u'_i 2^{(3 \log b)i}$ 和 $\hat{v} = \sum_{i=0}^{b-1} v'_i 2^{(3 \log b)i}$;

c) 利用 2.6 节中的算法计算积 $\hat{u}\hat{v}$;

d) 乘积 $\hat{u}\hat{v} = \sum_{i=0}^{2b-1} y'_i 2^{(3 \log b)i}$, 其中 $y'_i = \sum_{j=0}^{2b-1} u'_j v'_{i-j}$ 。 w_i 模 b 的值可以通过计算 $w'_i = (y'_i - y'_{b+i})$ 模 b 恢复, 其中 $0 \leq i < b$;

5. 用下面的公式准确计算 w_i , 其中 w_i 在 $(b-1-i)2^{2l}$ 和 $(i+1)2^{2l}$ 之间;

$$w_i = (2^{2l} + 1)((w'_i - w'_i) \bmod b) + w'_i$$

6. 计算 $\sum_{j=0}^{b-1} w_j 2^{2lj}$ 模 $(2^n + 1)$, 这是期望的结果。 \square

定理 7.7 算法 7.3 计算 uv 模 $(2^n + 1)$ 的值。

证明: 由定理 7.2, 算法 7.3 的第 1 步到第 3 步正确地计算了 w_i 模 $2^{2l} + 1$ 的值。作为练习, 读者可自行证明算法的第 4 步计算 w_i 模 b 的值, 第 5 步计算 w_i 模 $b(2^{2l} + 1)$, 即确切的 w_i 值。 \square

定理 7.8 算法 7.3 所需要的执行时间为 $O_b(n \log n \log \log n)$ 步

证明: 由定理 7.6 的推论可知, 第 1 步到第 3 步需要的时间为 $O_b[bl \log b + bM(2l)]$, 其中 $M(m)$ 是通过算法的递归应用计算两个 m 位整数相乘所需要的时间。在第 4 步, 构造长度为 $3b \log b$ 的 \hat{u} 和 \hat{v} , 并且在 $O_b[(3b \log b)^{1.59}]$ 步内计算出乘积。在 b 足够大时, $(3b \log b)^{1.59} < b^2$, 所以第 4 步所需要的时间可以在第 1 到第 3 步中的 $O_b(b^2 \log b)$ 项下忽略不计。第 5、6 步都只需要 $O(n)$ 时间, 所以也可以忽略。

因为 $n = bl$ 且 b 最大为 \sqrt{n} , 所以, 可以得到递推公式如下

$$M(n) \leq cn \log n + bM(2l) \quad (7-9)$$

对于常量 c 和足够大的 n , 假设 $M'(n) = M(n)/n$ 。那么式(7-9)变成

$$M'(n) \leq c \log n + 2M'(2l) \quad (7-10)$$

又因为 $l \leq 2\sqrt{n}$,

$$M'(n) \leq c \log n + 2M'(4\sqrt{n}) \quad (7-11)$$

对于恰当的 c' , 该式蕴涵着 $M'(n) \leq c' \log n \log \log n$ 。在这里, 在式(7-11)中把 $M'(4\sqrt{n})$ 置换为 $c' \log(4\sqrt{n}) \log \log 4\sqrt{n}$ 。直接相乘可得

$$M'(n) \leq c \log n + 4c' \log \left(2 + \frac{1}{2} \log n \right) + c' \log n \log \left(2 + \frac{1}{2} \log n \right)$$

对于大的 n , $2 + \frac{1}{2} \log n \leq \frac{2}{3} \log n$ 。由此

$$M'(n) \leq c \log n + 4c' \log \frac{2}{3} + 4c' \log \log n + c' \log \frac{2}{3} \log n + c' \log n \log \log n \quad (7-12)$$

对比较大的 n 和足够大的 c' , 前四项可被负的第四项消去。所以式(7-12)简化为 $M'(n) \leq c' \log n$

$\log \log n$ 。从这里可以总结出结论 $M(n) \leq c'n \log n \log \log n$ 。 □

习题

- 7.1 在复数环中, 当 $n=3, 4, 5$ 时, 其主 n 次单位根是什么?
- 7.2 不使用临时数组 S , 试描述算法 7.2 如何实现。
- 7.3 利用下面序列计算复数环下的离散傅里叶变换。
- a) $[0, 1, 2, 3]$
b) $[1, 2, 0, 2, 0, 0, 0, 1]^*$
- 7.4 当 n 不是 2 的幂时, 推广快速傅里叶变换算法。
- 定义** 如果 ω 和 n 有乘法的逆, ω 是模 m 整数环的主 n 次单位根时, 整数三元组 (ω, n, m) 称为是可接纳的 (admissible)。
- 7.5 下面哪个三元组是可接纳的。
- a) $(3, 4, 5)$
b) $(2, 6, 21)$
c) $(2, 6, 7)$
- 7.6 证明, 如果三元组 (ω, n, m) 是可接纳的, 那么当 $1 \leq p < n$ 时, $\omega^n \equiv 1$ 模 m 并且 $\omega^p \not\equiv 1$ 。
- *7.7 证明, 如果 n 是 2 的幂, $2^n \equiv 1$ 模 m , $2^p - 1$ 和 m 互质, 其中 $1 \leq p < n$, 那么 $(2, n, m)$ 是可接纳的三元组。
- **7.8 证明, 如果 m 是质数并且 ω 是任意数, 那么存在一个 n , 使 (ω, n, m) 是可接纳的三元组。
- *7.9 证明, 如果 a 和 b 互质, 那么存在 c 使 $ac \equiv 1$ 模 b 。同时逆命题也成立。证明模 b 下 c 是唯一的。
- 7.10 计算 $(10101110011110)_2$ 模 $2^5 + 1$ 。
- 7.11 假设 t 是以 10 为基的整数。证明将 t 的各位数字相加, 然后将结果的各位数字相加, 如此下去, 直到剩下一位数字, 此时结果等于 t 模 9。
- 7.12 使用卷积定理和可接纳三元组 $(2, 8, 17)$ 计算序列 $[1, 2, 3, 4]$ 和 $[4, 3, 2, 1]$ 模 17 的卷积。
- 7.13 利用算法 7.3 计算二进制数 $(1011011)_2$ 和 $(10001111)_2$ 的乘积。
- 7.14 证明: 对 $n=2^l$ 求平方根 $\log \log n$ 次得到的结果是 2。
- **7.15 利用卷积而不是包卷积设计一个快速整数乘法算法。分析算法的近似运行时间。
- *7.16 序列 $\mathbf{a} = [a_0, a_1, \dots, a_{n-1}]^T$ 的循环差分 (cyclic difference) 表示为 $\Delta \mathbf{a}$, 是序列 $[a_0 - a_{n-1}, a_1 - a_0, a_2 - a_1, \dots, a_{n-1} - a_{n-2}]^T$ 。假设 $F(\mathbf{a}) = [a'_0, a'_1, \dots, a'_{n-1}]^T$ 。证明 $F(\Delta \mathbf{a}) = [0, a'_1(1 - \omega), a'_2(1 - \omega^2), \dots, a'_{n-1}(1 - \omega^{n-1})]$, 其中 ω 是 n 次单位根。
- **7.17 利用习题 7.16 的结果证明, 如果 $X(n) - X(n-1) = n$ 和 $X(0) = 0$, 那么 $X(n) = n(n+1)/2$ 。
- *7.18 循环行列式是一个矩阵, 该矩阵的每一行都是上面一行元素向右位移一个单位。例如,
- $$\begin{bmatrix} a & b & c \\ c & a & b \\ b & c & a \end{bmatrix}$$
- 是一个 3×3 循环行列式。证明, 计算长度为 n 的向量的离散傅里叶变换等价于 $(n-1) \times (n-1)$ 循环行列式相乘, 其中 n 是质数。
- **7.19 证明有限域上的有限傅里叶变换可以在 $O_\lambda(n \log n)$ 步内计算得到, 其中 n 是质数。
- *7.20 思考通过每一点上的导数的值表示一个多项式, 是否从系数到导数的值变换是线性的?

- **7.21 根据多项式操作给卷积一个“物理”定义。
- *7.22 利用快速傅里叶变换编写一个算法在 $O(n \log n)$ 时间内计算出特普利茨矩阵和向量的乘积。把解决方法和习题 6.26(b) 相对比。

研究性练习题

- 7.23 寻找一个整数相乘或者离散傅里叶变换的快速算法。证明在一定模型的限制下, Schönhage-Strassen 算法或者快速傅里叶变换是最好的。在特定的约束下, Paterson, Fischer 和 Meyer [1974] 证明只要 $O_b((n \log n)/(\log \log n))$ 时间就可以实现按位整数乘法。类似, Morgenstern [1973] 说明了在特定的限制条件下, 离散傅里叶变换需要 $O_A(n \log n)$ 时间。

文献与注释

Cooley、Lewis 和 Welch [1967] 对快速傅里叶变换的研究是在 Runge 和 König [1924] 的基础上进行的。这个技术被 Danielson 和 Lanczos [1942] 和 Good [1958, 1960] 应用过。Cooley 和 Tukey [1965] 基础论文基本上清晰描述了该技术的本质。把 FFT 技术发扬用到多项式除法, 主要归功于 Fiduccia [1972]。由于在计算上的重要性, 大量的注意力吸引到有效实现的算法上。例如 Gentleman 和 Sande [1966] 以及 Rabiner 和 Rader [1972] 的许多文章都是有关这方面的。整数乘法算法由 Schönhage 和 Strassen [1971] 提出。习题 7.18 和 7.19 取自 Rader [1968]。

第8章 整数与多项式计算

把整数和多项式计算合在一起讲的主要原因在于大多数处理整数和一元多项式的算法在本质上是相同的。不仅是乘法和除法,还有更复杂的操作也很相似。例如,求一个整数模第二个整数的余数的运算等同于对多项式在某一点求值。利用余数表示一个整数的方法等同于用一些点的值来表示多项式。从余数构造一个整数的过程(“中国剩余定理”)等价于插值一个多项式。

本章将介绍一些特定的整数和多项式操作,例如与乘法操作具有相同时间复杂度要求的除法和开方运算。其他操作如上面提到的余数操作或计算最大公因子操作,这些操作需要至多比乘法操作多一个 $\log n$ 因子的时间,其中 n 是用二进制整数的长度或者多项式的次数。本章对整数和多项式的叙述是交替进行的,通常只证明其中一种情况,另外一种则留作练习。和其他章一样,重点在于提供与目前已知最好的相接近的算法。

本章最后将简要讨论大部分系数不为零的多项式(稠密多项式)和大多数系数为零的多项式(稀疏多项式)模型之间的区别。在处理拥有许多变量的多项式(本书中不讨论这种情况)时,稀疏模型非常有用。

8.1 整数和多项式的相似性

非负整数和一元多项式之间最明显的相似性在于它们都可以表示为有限的幂数列 $\sum_{i=0}^{n-1} a_i x^i$ 。对于整数而言,如果 $x=2$,那么 a_i 可以从集合 $\{0, 1\}$ 中取值。对于多项式而言,如果 x 不确定,那么 a_i 可以从某系数集^①中取值。

表示整数或多项式的幂数列长度基本上是自然的“规模”衡量尺度。对于二进制整数而言,尺度是表示该整数需要的位数;对于多项式而言,是多项式的系数的个数。有如下定义。

定义 如果 i 是一个非负整数,那么 $\text{SIZE}(i) = \lfloor \log i \rfloor + 1$ 。如果 $p(x)$ 是一个多项式,那么 $\text{SIZE}(p) = \text{DEG}(p) + 1$,其中 $\text{DEG}(p)$ 是 p 的次数,即系数非零的 x 的最高幂。

对于整数和多项式,可以进行类似的划分。如果 a 和 b 是两个整数,且 $b \neq 0$,那么存在唯一的一对整数 q 和 r ,满足

$$1. a = bq + r, \text{ 且}$$

$$2. r < b$$

其中 q 和 r 分别是 a 除以 b 的商和余数。

类似地,如果 a 和 b 是两个多项式,且 b 不是常量,那么存在唯一的 q 和 r 多项式满足

$$1. a = bq + r, \text{ 且}$$

$$2. \text{DEG}(r) < \text{DEG}(b)。$$

整数和多项式之间另一个相似性是它们都存在非常快速的乘法算法。上一章讨论了两个具有实数系数的 n 次多项式可以通过使用 FFT 在 $O_A(n \log n)$ 时间内实现乘法操作。因为在实际应用中,使用多项式的固定精度的系数表示多项式,并且通过系数的算术操作来实现各种多项式操作,所以算术操作的复杂性计算是合理的。

① 在后面,可以认为系数集是实数域(见 12.1 节),由于算法复杂性的计算是根据系数域上的操作次数而定的,本章的结果可以应用于任何域。由此,不必考虑模型中的系数大小。

如果使用 7.5 节中的 Schönhage-Strassen 算法, 那么两个 n 位整数乘法可以在 $O_B(n \log n \log \log n)$ 时间内完成。对于整数而言, 位操作的次数是考虑算法复杂度时唯一感兴趣的量度。在以下两种情况下都不把整数乘法当作原语操作。第一, 在设计乘法操作硬件时, 位操作次数反映一个乘法周期需要的元素数量; 第二, 在为定长字计算机设计定点的任意精度算法时。如在进行 n 位精度乘法时, 位操作次数和所需的机器指令数目是相关的。

如此, 对于所用的两种不同的复杂度(算术和位), 多项式和整数算术计算的结果将显示出很高的相似性。进一步讲, 用幂数列表示整数时, 位操作就是系数操作, 对多项式而言, 算术操作也是系数操作, 两种量度是很相似的。

8.2 整数的乘法和除法

下面讨论位操作的整数乘法需要的时间和整数除法需要的时间只差一个常数因子, 近似于整数开方和求倒数需要的时间。本节将用到的四个符号及其解释如下:

符号	含义
$M(n)$	长度为 n 的两个整数乘法需要的时间
$D(n)$	长度小于 $2n$ 的整数除以长度为 n 的整数需要的时间
$S(n)$	长度为 n 的整数开方需要的时间
$R(n)$	长度为 n 的整数求倒数需要的时间

每一种情况下, 需要的时间以位操作为单位。同时合理地假设 $M(n)$ 满足如下条件:

$$a^2 M(n) \geq M(an) \geq aM(n)$$

其中 $a \geq 1$ 。其他三个函数也满足。

首先证明 n 位整数 i 的倒数可以在两个 n 位整数相乘需要的时间内计算完。因为 $1/i$ 在 $i > 1$ 时的结果不是一个整数, 所以求 i 的“倒数”的真正含义是求逼近分数 $1/i$ 的 n 位有效值。由于假设比例的改变(二进制小数点的位移)不花费时间, 因此可以说 i 的“倒数”等同于 2^{2n-1} 除以 i 的商。在本节的剩余部分用项的倒数(term reciprocal)表示这种含义。

首先考虑为数 $A = 1/P$ 寻找一个逼近序列, 其中 P 是个整数。假设 A_i 是对 $1/P$ 的第 i 个逼近。那么 A 的确切值可以表示为

$$A = A_i + (1/P)(1 - A_i P) \quad (8-1)$$

如果用式(8-1)中的 A_i 逼近 $1/P$, 可以得到

$$A_{i+1} = A_i + A_i(1 - A_i P) = 2A_i - A_i^2 P \quad (8-2)$$

该迭代公式可以根据第 i 个逼近计算出第 $(i+1)$ 个逼近。注意, 如果 $A_i P = 1 - S$, 那么

$$A_{i+1} P = 2A_i P - A_i^2 P^2 = 2(1 - S) - (1 - S)^2 = 1 - S^2$$

这说明了迭代公式(8-2)是二倍收敛的(converge quadratically)。如果 $S \leq \frac{1}{2}$, 那么每次迭代得到的正确位数是双倍增长的。

由于 P 可能有许多位, 所以没有必要把 P 的所有位用在早期的逼近中。因为只有 P 的前导位(leading bit)会影响正确计算正确的 A_{i+1} 位。进一步, 如果 A_{i+1} 到小数点右边的头 k 位是正确的, 那么利用公式(8-2)可以得到 A_{i+1} 的头 $2k$ 位。也就是说计算 $A_{i+1} = 2A_i - A_i^2 P$ 时, $2A_i$ 截取小数点右边的 k 位, $A_i^2 P$ 取 P 的 $2k$ 位, 然后对得到的结果在小数点右边 $2k$ 位置截断。因为在前面假设 P 是确切的, 所以最后得到的近似值肯定影响收敛性。

下面的算法中将用到这个思想。实际计算 $\lfloor 2^{2n-1}/P \rfloor$ 的商, 其中 $P = p_1 p_2 \cdots p_n$ 是一个 n 位整数, 且 $p_1 = 1$ 。使用的方法本质上就是公式(8-2), 引入比例(scaling)能够对整数进行专门处理。

因此,没有必要保持小数点。

算法 8.1 求整数的倒数。

输入: n 位整数 $P = [p_1 p_2 \cdots p_n]$, 其中 $p_1 = 1$ 。为了叙述方便, 假设 n 是 2 的幂, 并且 $[x]$ 是位串 x 所表示的整数(例如 $[110] = 6$)。

输出: 整数 $A = [a_0 a_1 \cdots a_n]$, 满足 $A = \lfloor 2^{2n-1}/P \rfloor$ 。

方法: 调用 $\text{RECIPROCAL}([p_1 p_2 \cdots p_n])$, 其中 RECIPROCAL 是图 8-1 中的递归过程。对于 2 的任意次幂 k , 该过程计算 $\lfloor 2^{2k-1}/[p_1 p_2 \cdots p_k] \rfloor$ 的逼近值。除了 P 是 2 的幂的结果是 $k+1$ 位整数外, 其他结果都是一个普通的 k 位整数。

给定一个 k 位 $[p_1 p_2 \cdots p_k]$ 的倒数, 第 2~4 行计算 $2k-3$ 位 $[p_1 p_2 \cdots p_{2k}]$ 的倒数。第 5~7 行校正最后三位。实际上可以跳过 5~7 行, 最后通过对公式(8-2)的应用可以得到想要的准确性。为了简化对算法的理解, 并且证明算法确实有效, 选择包含第 5~7 行的循环。□

例 8.1 计算 $2^{15}/153$ 。这里用

$$n = 8 \text{ 和 } 153 = [p_1 p_2 \cdots p_8] = [10011001]$$

调用

$\text{RECIPROCAL}([10011001])$

递归地按顺序调用 RECIPROCAL , 分别使用参数 $[1001]$, $[10]$ 和 $[1]$ 。在第 1 行得到 $\text{RECIPROCAL}([1]) = [10]$, 在第 2 行返回 $\text{RECIPROCAL}([10])$, 其间设置 $[c_0 c_1] \leftarrow [10]$ 。然后在第 3 行, 计算 $[d_1 \cdots d_4] \leftarrow [10] * 2^3 - [10]^2 * [10] = [1000]$ 。接下去在第 4 行设 $[a_0 a_1 a_2] = [100]$ 。在 5~7 行的循环中没有变化。返回到 $\text{RECIPROCAL}([1001])$, 其中 $[100]$ 是 $2^3/[10]$ 的近似值。

返回第 2 行, 当 $k=4$ 时, 有 $[c_0 c_1 c_2] = [100]$ 。那么 $[d_1 \cdots d_8] = [01110000]$ 且 $[a_0 \cdots a_4] = [01110]$ 。这次在第 5~7 行的循环中仍旧没有变化。返回第 2 行的 $\text{RECIPROCAL}([10011001])$, 且 $[c_0 \cdots c_4] = [01110]$, 那么

$$[d_1 \cdots d_{16}] = [0110101011011100]$$

这样在第 4 行, $[a_0 \cdots a_8] = [011010101]$ 。在第 6 行, 发现 $[011010101] * [10011001]$ 是十进制的 $213 * 153$, 它的值为 32589, 而 $2^{15} = 32768$ 。这样在第 5~7 行的循环中, 1 被加到 213 上, 得到 214, 或者是二进制的 $[011010110]$ 。□

```

procedure RECIPROCAL( $[p_1 p_2 \cdots p_k]$ ):
1.  if  $k = 1$  then return  $[10]$ 
    else
        begin
2.       $[c_0 c_1 \cdots c_{k/2}] \leftarrow \text{RECIPROCAL}([p_1 p_2 \cdots p_{k/2}]);$ 
3.       $[d_1 d_2 \cdots d_{2k}] \leftarrow [c_0 c_1 \cdots c_{k/2}] * 2^{2k/2} -$ 
          $[c_0 c_1 \cdots c_{k/2}]^2 * [p_1 p_2 \cdots p_k];$ 
         comment 虽然右边第3行看起来产生一个  $(2k+1)$  位数,
           实际上头  $(2k+1)$  位总是0;
4.       $[a_0 a_1 \cdots a_k] \leftarrow [d_1 d_2 \cdots d_{k+1}];$ 
         comment  $[a_0 a_1 \cdots a_k]$  是  $2^{2k-1}/[p_1 p_2 \cdots p_k]$  非常好的近
           似值。下面的循环在必要的情况下通过加入最后三位
           提高逼近程度。
5.      for  $i \leftarrow 2$  step  $-1$  until 0 do
6.          if  $([a_0 a_1 \cdots a_k] + 2^i) * [p_1 p_2 \cdots p_k] \leq 2^{2k-1}$  then
7.               $[a_0 a_1 \cdots a_k] \leftarrow [a_0 a_1 \cdots a_k] + 2^i;$ 
8.      return  $[a_0 a_1 \cdots a_k]$ 
        end

```

图 8-1 计算整数倒数的过程

定理 8.1 算法 8.1 得到 $[a_0 a_1 \cdots a_k]$, 满足

$$[a_0 a_1 \cdots a_k] * [p_1 p_2 \cdots p_k] = 2^{2k-1} - S$$

且 $0 \leq S < [p_1 p_2 \cdots p_k]$ 。

证明: 对 k 进行归纳。归纳基础 $k=1$, 由第 1 行, 容易得证。对于归纳步, 假设 $C = [c_0 c_1 \cdots c_{k/2}]$, $P_1 = [p_1 p_2 \cdots p_{k/2}]$, $P_2 = [p_{k/2+1} p_{k/2+2} \cdots p_k]$ 。那么 $P = [p_1 p_2 \cdots p_k] = P_1 2^{k/2} + P_2$ 。由递归假设

$$CP_1 = 2^{k-1} - S$$

其中 $0 \leq S < P_1$ 。在第 3 行, $D = [d_1 d_2 \cdots d_{2k}]$ 由下式得到

$$D = C 2^{3k/2} - C^2 (P_1 2^{k/2} + P_2) \quad (8-3)$$

因为 $p_1 = 1$, $P_1 \geq 2^{k/2-1}$, 所以有 $C \leq 2^{k/2}$ 。由于 $D < 2^{2k}$, 所以 $2k$ 位足以用来表示 D 。

考虑乘积 $PD = (P_1 2^{k/2} + P_2)D$, 由式(8-3)有:

$$PD = CP_1 2^{2k} + CP_2 2^{3k/2} - (CP_1 2^{k/2} + CP_2)^2 \quad (8-4)$$

通过把式(8-4)中的 CP_1 替换为 $2^{k-1} - S$, 并且进行一些代数化简, 可得:

$$PD = 2^{3k-2} - (S 2^{k/2} - CP_2)^2 \quad (8-5)$$

式(8-5)除以 2^{k-1} , 可得

$$2^{2k-1} = \frac{PD}{2^{k-1}} + T \quad (8-6)$$

其中 $T = (S 2^{k/2} - CP_2)^2 2^{-(k-1)}$ 。通过归纳假设和 $P_1 < 2^{k/2}$, 可得 $S < 2^{k/2}$ 。因为 $C \leq 2^{k/2}$ 和 $P_2 < 2^{k/2}$, 有 $0 \leq T < 2^{k+1}$ 。

在第 4 行, $A = [a_0 a_1 \cdots a_k] = \lfloor D/2^{k-1} \rfloor$ 。现在有

$$P \lfloor \frac{D}{2^{k-1}} \rfloor > P \left(\frac{D}{2^{k-1}} - 1 \right)$$

所以, 从式(8-6)可得

$$2^{2k-1} \geq \frac{PD}{2^{k-1}} \geq P \lfloor \frac{D}{2^{k-1}} \rfloor > \frac{PD}{2^{k-1}} - P = 2^{2k-1} - T - P > 2^{2k-1} - 2^{k+1} - 2^k$$

这样

$$P \lfloor \frac{D}{2^{k-1}} \rfloor = 2^{2k-1} - S'$$

其中 $0 \leq S' < 2^{k+1} + 2^k$ 。因为 $P \geq 2^{k-1}$, 至多把 6 加到 $\lfloor D/2^{k-1} \rfloor$ 上, 可得到满足归纳假设 k 的数。因为这部分工作被第 5~7 行执行, 所以归纳步成立。□

定理 8.2 存在一个常数 c , 使 $R(n) \leq cM(n)$ 成立。

证明: 算法 8.1 可以在 $O_b(M(n))$ 时间内完成。第 2 行要求 $R(k/2)$ 次位操作。第 3 行包含平方和乘法操作, 分别需要 $M(k/2+1)$ 和 $M(k+2)$ 时间, 再加上一个减法需要 $O_b(k)$ 时间。注意, 乘以 2 的幂的乘法操作不需要任何位操作; 乘数的位数可以简单认为是占有新的位置, 即好像已经位移了。由对 M 的假设, $M(k/2+1) \leq \frac{1}{2}M(k+2)$ 成立。进一步而言, $M(k+2) - M(k)$

的时间复杂度即 $O_b(k)$ (参见 2.6 节为例), 那么存在常数 c' , 使第 3 行的执行时间界为 $\frac{3}{2}M(k) + c'k$ 。第 4 行的执行时间显然为 $O_b(k)$ 。

在第 5~7 行的循环中需要三个乘法, 但是这些计算可以通过一个乘法 $[a_0 a_1 \cdots a_k] * [p_1 p_2 \cdots p_k]$ 和最多 $2k$ 位整数的加法和减法来完成。所以第 5~7 行的执行时间界为 $M(k) + c''(k)$, c'' 是一个常量。把所有代价相加, 有下式

$$R(k) \leq R\left(\frac{k}{2}\right) + \frac{5}{2}M(k) + c_1k \quad (8-7)$$

其中 c_1 是一常量。

可以说存在一个常量 c , 使得 $R(k) \leq cM(k)$ 成立。选择 c , 满足 $c \geq R(1)/M(1)$ 和 $c \geq 5 + 2c_1$ 。对 k 进行归纳来验证这一说法。初始 $k=1$, 结论显然成立。对式(8-7)进行归纳, 因为

$$R(k) \leq cM\left(\frac{k}{2}\right) + \frac{5}{2}M(k) + c_1k \quad (8-8)$$

因为对 M 的假设可以得到 $M(k/2) \leq \frac{1}{2}M(k)$, 并且显然 $k \leq M(k)$, 可以把式(8-8)改写为

$$R(k) \leq \left(\frac{c}{2} + \frac{5}{2} + c_1\right)M(k) \quad (8-9)$$

因为 $c \geq 5 + 2c_1$, 所以式(8-9)隐含 $R(k) \leq cM(k)$ 。□

显然, 如果 P 的位数不考虑, 二进制小数点的具体位置, 算法 8.2 可以用来计算 $1/P$ 到 n 个有效位。例如, 如果 $\frac{1}{2} < P < 1$, P 有 n 位, 很显然, 通过改变比例, 可使 $1/P$ 结果为 1 加上小数点后面的 $n-1$ 个有效位。

下面说明对一个大小为 n 的整数进行平方操作需要的时间 $S(n)$ 不会比求一个大小为 n 的整数的倒数需要的时间 $R(n)$ 大一个数量级。该技术用到恒等式

$$P^2 = \frac{1}{\frac{1}{P} - \frac{1}{P+1}} - P \quad (8-10)$$

下面的算法在恰当的比例下使用式(8-10)。

算法 8.2 通过倒数求平方。

输入: 以二进制表示的 n 位整数 P 。

输出: P^2 的二进制表示。

方法:

1. 追加 $2n$ 个 0 到 P , 再利用算法 8.1 计算 $A = \lfloor 2^{4n-1}/P \rfloor$ 。然后利用 RECIPROCAL[⊙] 计算 $\lfloor 2^{6n-1}/P2^{2n} \rfloor$, 再移位;
2. 用类似的方法计算 $B = \lfloor 2^{4n-1}/(P+1) \rfloor$;
3. 假设 $C = A - B$ 。 $C = 2^{4n-1}/(P^2 + P) + T$, 其中 $|T| \leq 1$ 。引入 T 在于, 当计算 A 和 B 时, 截取操作将会引起大于 1 的误差。因为 $2^{2n-2} < P^2 + P < 2^{2n}$, 所以 $2^{2n+1} \geq C \geq 2^{2n-1}$;
4. 计算 $D = \lfloor 2^{4n-1}/C \rfloor - P$;
5. 设 Q 是 P 的最后 4 位, 尽可能小地向上或向下调整 D 的最后 4 位, 使得值能够和 Q^2 的最后 4 位相一致。□

定理 8.3 算法 8.2 计算 P^2 。

证明: 由于第 1 和 2 步中的截取操作, 我们能够确定的仅是

$$C = \frac{2^{4n-1}}{P^2 + P} + T, \text{ 其中 } |T| \leq 1$$

由于 $2^{2n-1} \leq C \leq 2^{2n+1}$, 并且 C 的误差在 $-1 \sim +1$ 范围内, 所以 $2^{4n-1}/C$ 的误差最多为

$$\left| \frac{2^{4n-1}}{C} - \frac{2^{4n-1}}{C-1} \right| = \left| \frac{2^{4n-1}}{C^2 - C} \right|$$

⊙ 在算法 8.1 中定义的 RECIPROCAL 是专门针对长度为 2 的幂的整数。对于长度不是 2 的幂的整数推广是在必要的时候通过加 0 和改变比例。

因为 $C^2 - C \geq 2^{4n-3}$, 所以误差最多为 4。第 4 行的截取操作会把误差扩大到 5。这样, $|P^2 - D| \leq 5$ 。所以在第 5 步中计算 P^2 的最后 4 位可以确保 D 调整到正好等于 P^2 。□

例 8.2 假设 $n=4$, $P=[1101]$ 。那么

$$A = \lfloor 2^{15} / [1101] \rfloor = [100111011000]$$

并且

$$B = \lfloor 2^{15} / [1110] \rfloor = [100100100100]$$

下面计算 C ,

$$C = A - B = [10110100]$$

然后,

$$D = \lfloor 2^{15} / C \rfloor - P = [10110110] - [1101] = [10101001]$$

因此, 13 的平方 $D=169$, 在第 5 步中不需要进行修正。□

定理 8.4 存在一个常数 c , 使 $S(n) \leq cR(n)$ 。

证明: 算法 8.2 使用三个长度最多为 $3n$ 的二进制串的倒数, 第 3 和 4 步的减法需要 $O_b(n)$ 时间, 加上第 5 步中与 n 无关的固定量的计算开销。所以

$$S(n) \leq 3R(3n) + c_1 n \quad (8-11)$$

其中 c_1 是常量。由此, $S(n) \leq 27R(n) + c_1 n$ 。因为 $R(n) \geq n$, 选择 $c = 27 + c_1$, 定理得证。□

定理 8.5 $M(n)$, $R(n)$, $D(n)$ 和 $S(n)$ 都是常量因子相关的。

证明: 已经证明了存在常量 c_1 和 c_2 , 使 $R(n) \leq c_1 M(n)$ 和 $S(n) \leq c_2 R(n)$ 成立。通过 $AB = \frac{1}{2} [(A+B)^2 - A^2 - B^2]$, 可以很容易地得到 $M(n) \leq c_3 S(n)$ 。由此, M 、 R 和 S 都是常量因子相关的。

当讨论 n 位数的除法时, 是指一个最多 $2n$ 位的数除以一个 n 位数, 得到最多 $n+1$ 位的结果。显然有 $R(n) \leq D(n)$, 所以 $M(n) \leq c_2 c_3 D(n)$ 。进一步通过恒等式 $A/B = A * (1/B)$, 可以知道存在常数 c , 使得(注意比例)

$$D(n) \leq M(2n) + R(2n) + cn \quad (8-12)$$

因为很容易可以得到 $R(2n) \leq c_1 M(2n)$, $M(2n) \leq 4M(n)$, 可改写式(8-12)为

$$D(n) \leq 4(1 + c_1)M(n) + cn \quad (8-13)$$

因为 $M(n) \geq n$, 从式(8-13)可以得到 $D(n) \leq c_4 M(n)$, 其中 $c_4 = 4 + 4c_1 + c$ 。因此, 所有的函数都是在 $dM(n)$ 和 $eM(n)$ 之间, 其中 d 和 e 是正整数常量。□

推论 $2n$ 位整数除以 n 位整数可以在 $O_b(n \log n \log \log n)$ 时间内完成。

证明: 通过运用定理 7.8 和 8.5 证明。□

8.3 多项式的乘法和除法

前一节所介绍的技术都是一元多项式算术。在本节中, $M(n)$, $D(n)$, $R(n)$ 和 $S(n)$ 分别表示 n 次多项式乘、除、求倒数和求平方所需要的时间。和前面一样, 假设对于 $a \geq 1$, 不等式 $a^2 M(n) \geq M(an) \geq aM(n)$ 成立, 其他函数类似。

对于 n 次多项式 $p(x)$, “倒数”表示为 $\lfloor x^{2n}/p(x) \rfloor^\ominus$ 。 $D(n)$ 为得到 $\lfloor s(x)/p(x) \rfloor$ 所耗费的时间, 其中 $p(x)$ 的次数为 n , $s(x)$ 的次数不超过 $2n$ 。注意, 可以通过乘以和除以 x 的幂来“伸缩”(scale)多项式, 和上一节通过 2 的幂伸缩整数一样。

⊖ 和整数中的概念类似, 可用“下取整函数”(floor function)表示多项式的商。也就是说, 如果 $p(x)$ 不是一个常量, 那么 $\lfloor s(x)/p(x) \rfloor$ 是唯一的 $q(x)$, 满足 $s(x) = p(x)q(x) + r(x)$ 和 $\text{DEG}(r(x)) < \text{DEG}(p(x))$ 。

由于本节中的许多结果和上节中有关整数的结果类似, 仅给出一个详细的结果: 多项式“倒数”算法, 该算法类似于算法 8.1 中整数的倒数算法。由于在多项式算法中不会像整数算法中那样产生把次数序列从一个地方移动到另外一个地方的情况, 所以多项式算法在某种程度上比整数算法更简单。这样, 多项式算法不需要和算法 8.1 中的第 5~7 行一样调整最低有效位的位置。

算法 8.3 多项式倒数。

输入: $n-1$ 次多项式 $p(x)$, 其中 n 是 2 的幂[即 $p(x)$ 存在 2^t 个项, 其中 t 是整数]。

输出: “倒数” $\lfloor x^{2n-2}/p(x) \rfloor$ 。

方法: 图 8-2 中定义一个新的过程

$$\text{RECIPROCAL} \left(\sum_{i=0}^{k-1} a_i x^i \right)$$

其中 k 是 2 的幂且 $a_{k-1} \neq 0$ 。该过程计算

$$\left\lfloor \frac{x^{2k-2}}{\sum_{i=0}^{k-1} a_i x^i} \right\rfloor$$

注意, 如果 $k=1$, 那么参数是一个常量 a_0 , 其倒数 $1/a_0$ 是另外一个常量。假设对系数的每个操作都可以在一步内完成, 且计算 $1/a_0$ 不需要调用 RECIPROCAL。

算法本身以 $p(x)$ 为参数调用 RECIPROCAL。

□

```

procedure RECIPROCAL  $\left( \sum_{i=0}^{k-1} a_i x^i \right)$ :
1.  if  $k = 1$  then return  $1/a_0$ 
   else
     begin
2.       $q(x) \leftarrow \text{RECIPROCAL} \left( \sum_{i=k/2}^{k-1} a_i x^{i-k/2} \right)$ ;
3.       $r(x) \leftarrow 2q(x)x^{(3/2)k-2} - (q(x))^2 \left( \sum_{i=0}^{k-1} a_i x^i \right)$ ;
4.      return  $\lfloor r(x)/x^{k-2} \rfloor$ 
   end

```

图 8-2 计算多项式倒数的算法

例 8.3 计算 $\lfloor x^{14}/p(x) \rfloor$, 其中 $p(x) = x^7 - x^6 + x^5 + 2x^4 - x^3 - 3x^2 + x + 4$ 。算法第 2 行计算 $x^3 - x^2 + x + 2$ 的倒数, 也就是 $q(x) = \lfloor x^6/(x^3 - x^2 + x + 2) \rfloor$ 。可以验证 $q(x) = x^3 + x^2 - 3$ 。因为 $k=8$, 算法第 3 行计算 $r(x) = 2q(x)x^{10} - (q(x))^2 p(x) = x^{13} + x^{12} - 3x^{10} - 4x^9 + 3x^8 + 15x^7 + 12x^6 - 42x^5 - 34x^4 + 39x^3 + 51x^2 - 9x - 36$ 。在第 4 行, 结果是 $s(x) = x^7 + x^6 - 3x^4 - 4x^3 + 3x^2 + 15x + 12$ 。记 $s(x)p(x)$ 为 x^{14} 加上一个 6 次多项式。

□

定理 8.6 算法 8.3 正确地计算多项式的倒数。

证明: 对 k 进行归纳, 其中 k 是 2 的幂。如果 $s(x) = \text{RECIPROCAL}(p(x))$, $p(x)$ 为 $k-1$ 次, 那么 $s(x)p(x) = x^{2k-2} + t(x)$, 其中 $t(x)$ 是次数小于 $k-1$ 的多项式。归纳基础 $k=1$ 时, 因为 $p(x) = a_0$, $s(x) = 1/a_0$, $t(x)$ 不需要, 所以结论显然。

在归纳步, 假设 $p(x) = p_1(x)x^{k/2} + p_2(x)$, 其中 $\text{DEG}(p_1) = k/2 - 1$, $\text{DEG}(p_2) \leq k/2 - 1$ 。那么由归纳假设, 如果 $s_1(x) = \text{RECIPROCAL}(p_1(x))$, 那么, $s_1(x)p_1(x) = x^{k-2} + t_1(x)$, 其中, $\text{DEG}(t_1) < k/2 - 1$ 。在第 3 行中计算

$$r(x) = 2s_1(x)x^{(3/2)k-2} - (s_1(x))^2(p_1(x)x^{k/2} + p_2(x)) \quad (8-14)$$

这足以说明 $r(x)p(x)$ 是 x^{3k-4} 加上次数小于 $2k-3$ 的项。那么在第 4 行中除以 x^{k-2} 可以得到期望的结果。

由式(8-14)和 $p(x) = p_1(x)x^{k/2} + p_2(x)$, 可以得到

$$r(x)p(x) = 2s_1(x)p_1(x)x^{2k-2} + 2s_1(x)p_2(x)x^{(3/2)k-2} - (s_1(x)p_1(x)x^{k/2} + s_1(x)p_2(x))^2 \quad (8-15)$$

在式(8-15)中把 $s_1(x)p_1(x)$ 替换为 $x^{k-2} + t_1(x)$, 可以得到

$$r(x)p(x) = x^{3k-4} - (t_1(x)x^{k/2} + s_1(x)p_2(x))^2 \quad (8-16)$$

因为 $\text{DEG}(t_1) < k/2 - 1$, $s_1(x)$ 和 $p_2(x)$ 的次数最大为 $k/2 - 1$, 所以式(8-16)中除了 x^{3k-4} 外, 其余项的次数最大为 $2k-4$ 。□

当考虑不同的时间复杂度衡量标准(分别为位操作和算术操作)时, 很明显算法 8.1 和 8.3 的运行时间是相似的。通过类似的方式, 把算术操作步骤替换为位操作, 可以得到 8.2 节中的时间边界也适用于多项式情况。由此, 有下面的定理。

定理 8.7 假设 $M(n)$, $D(n)$, $R(n)$ 和 $S(n)$ 分别表示一元多项式乘法、除法、求倒数和平方的算术复杂度, 那么这些函数是常数相关的。

证明: 过程类似于定理 8.5, 得到的结论也与之相同。□

推论: $2n$ 次多项式除以 n 次多项式可以在 $O_A(n \log n)$ 时间内完成。

证明: 由定理 7.4 的推论 3 和定理 8.7 可以得到。□

8.4 模算术

许多应用中用“模”的记法(“modular” notation)做整数的算术操作是非常便捷的。也就是说, 不用固定基数记法表示整数, 而是用它模一系列两两互质的整数得到的余数来表示。如果 p_0, p_1, \dots, p_{k-1} 是两两互质的整数, 并且 $p = \prod_{i=0}^{k-1} p_i$, 那么可以把任意 $u (0 \leq u < p)$ 表示为余数集合 u_0, u_1, \dots, u_{k-1} , 其中 $u_i = u \bmod p_i (0 \leq i < k)$ 。已知 p_0, p_1, \dots, p_{k-1} 序列时, 记 $u \leftrightarrow (u_0, u_1, \dots, u_{k-1})$ 。

对于结果在 $0 \sim p-1$ 之间的加、减、乘演算相当容易(另外, 也可以认为这些演算是在模 p 下进行的)。也就是说

$$u \leftrightarrow (u_0, u_1, \dots, u_{k-1}) \text{ 和 } v \leftrightarrow (v_0, v_1, \dots, v_{k-1})$$

那么

$$u + v \leftrightarrow (w_0, w_1, \dots, w_{k-1}), \text{ 其中 } w_i = (u_i + v_i) \bmod p_i \quad (8-17)$$

$$u - v \leftrightarrow (x_0, x_1, \dots, x_{k-1}), \text{ 其中 } x_i = (u_i - v_i) \bmod p_i \quad (8-18)$$

$$uv \leftrightarrow (y_0, y_1, \dots, y_{k-1}), \text{ 其中 } y_i = u_i v_i \bmod p_i \quad (8-19)$$

例 8.4 假设 $p_0 = 5, p_1 = 3, p_2 = 2$ 。因为 $4 = 4 \bmod 5, 1 = 4 \bmod 3, 0 = 4 \bmod 2$, 所以有 $4 \leftrightarrow (4, 1, 0)$ 。类似地, $7 \leftrightarrow (2, 1, 1)$ 和 $28 \leftrightarrow (3, 1, 0)$ 。通过上面公式(8-19)可以得到 $4 \times 7 \leftrightarrow (3, 1, 0)$, 是 28 的对应表示。也就是说 4×7 的第一个分量是 $4 \times 2 \bmod 5$ 得到 3; 第二个分量是 $1 \times 1 \bmod 3$ 得到 1; 最后一个分量是 $0 \times 1 \bmod 2$ 得到 0。类似地, $4 + 7 \leftrightarrow (1, 2, 1)$, 是 11 的对应表示; $7 - 4 \leftrightarrow (3, 0, 1)$ 是 3 的对应表示。□

然而, 我们还没有说明如何利用模算术快捷地计算除法操作。需要注意 u/v 不一定是整数, 即使是, 也不能够通过对每个 i 计算 $(u_i/v_i) \bmod p_i$ 找到其模的表示。实际上, 如果 p_i 不是质数, 那么存在多个 0 和 $p_i - 1$ 之间的 w , 其中 w 等于 $(u_i/v_i) \bmod p_i$, 因为 $wv_i \equiv u_i \bmod p_i$ 。例如, 如果 $p_i = 6, v_i = 3, u_i = 3$, 因为 $1 \times 3 \equiv 3 \times 3 \equiv 5 \times 3 \equiv 3 \bmod 6$, 所以 w 可以是 1、3 或者 5。在这种情况下, $(u_i/v_i) \bmod p_i$ 就变得“不合理”。

模表示方式主要的优点在于其算术操作所需要的硬件支持比卷积操作需要的硬件支持更少。其原因在于每个模操作的演算是不相关的。通常的数值表示(即基于基位的表示法)是不需要进位的。遗憾的是,它在除法操作和有效地发现溢出(也就是分辨出结果是否超出 $0 \sim p-1$ 这个范围)的问题难以克服,同时这类系统通常很少以计算机硬件实现。

尽管如此,这个思想还是有用的,主要是在多项式领域。在这里,我们希望找到一个不需要多项式除法的情况。下一节将会看到多项式计值和多项式余数(模其他多项式)的计算是紧密相关的。首先证明整数取模的算术操作会按预期进行。

证明的第一部分说明公式(8-17)、(8-18)和(8-19)成立。这些关系是显而易见的,所以留给读者作为练习。证明的第二部分说明 $u \leftrightarrow (u_0, u_1, \dots, u_{k-1})$ 这种对应关系是一对一的(同构)。虽然这个结论并不难证明,我们还是给出一个引理。

引理 8.1 假设 p_0, p_1, \dots, p_{k-1} 是一个整数集合,其中任何两个数都是互质的。设

$$p = \prod_{i=0}^{k-1} p_i$$

并且 $u_i = u \text{ 模 } p_i$ 。那么对于 $0 \leq u < p$, $u \leftrightarrow (u_0, u_1, \dots, u_{k-1})$ 和整数 u 是一对一关系,并且

$$(u_0, u_1, \dots, u_{k-1}), 0 \leq u_i < p_i$$

其中 $0 \leq i < k$ 。

证明:很显然,每个 u 存在一个对应的 k 元组。因为区间中正好有 u 的 p 种值,也正好有 p 个 k 元组,这足以证明每个 k 元组至多有一个整数 u 与之对应。假设 u, v 满足 $0 \leq u < v < p$,这两个数都和 k 元组 u_0, u_1, \dots, u_{k-1} 对应。那么 $v-u$ 必定是每个 p_i 的倍数。因为 p_i 是互质的, $v-u$ 必定是 p 的倍数。因为 $u \neq v$ 并且 $v-u$ 是 p 的倍数,那么 u 和 v 之间的差距至少是 p ,所以,它们不可能都在 $0 \sim p-1$ 范围内。□

为了使用关于模运算,需要有算法能够实现基数表示法和模表示法之间的转换。把整数 u 从基数表示法转换到模表示法的一个方法是用每个 p_i 除 u ,其中 $0 \leq i < k$ 。

假设在二进制表示法中,每个 p_i 需要 b 位,那么

$$p = \prod_{i=0}^{k-1} p_i$$

大致需要 bk 位,用 u 除以每个 p_i ,其中 $0 \leq u < p$,需要调用 k 次 kb 位整数除以 b 位整数的除法。把每个除法分为 k 个 $2b$ 位整数除以 b 位整数的方式,可以在 $O_b(k^2 D(b))$ 时间内实现向模表示方式的转换,其中 $D(n)$ 是整数除法需要的时间[由定理 8.5 的推论知至多为 $O_b(n \log n \log \log n)$]。

然而,回想 7.2 节中用到的多项式除法技术,利用该方法我们可以在更少的时间内完成该工作。和用整数 u 模 p_0, p_1, \dots, p_{k-1} 的 k 个除法不同,首先计算 $p_0 p_1, p_2 p_3, \dots, p_{k-2} p_{k-1}$ 的乘积,然后计算 $p_0 p_1 p_2 p_3, p_4 p_5 p_6 p_7, \dots$ 乘积,等等。下一步用分治法计算余数。通过用 u 模 $p_0 \dots p_{k/2-1}$ 和 u 模 $p_{k/2} \dots p_{k-1}$ 的除法分别得到余数 u_1 和 u_2 。计算 u 模 p_i ($0 \leq i < k$)的问题现在简化为两个规模减半的问题了。也就是说,当 $0 \leq i < k/2$ 时, $u \text{ 模 } p_i = u_1 \text{ 模 } p_i$,当 $k/2 \leq i < k$ 时, $u \text{ 模 } p_i = u_2 \text{ 模 } p_i$ 。

算法 8.4 余数的计算。

输入:模 P_0, p_1, \dots, p_{k-1} 和整数 u ,其中 $0 \leq u < p = \prod_{i=0}^{k-1} p_i$ 。

输出: u_i ($0 \leq i < k$),其中 $u_i = u \text{ 模 } p_i$ 。

方法:假设 k 是 2 的幂,也就是 $k=2^t$ 。(如果需要,可以加入一些 1 到输入作为模,使得 k 成为 2 的幂。)从计算模的乘积开始,方法类似于 7.2 节中 q_{lm} 的计算。

对于 $0 \leq j < t$, i 是 2^j 的倍数,且 $0 \leq i < k$, 设

$$q_{ij} = \prod_{m=i}^{i+2^j-1} p_m$$

那么 $q_{i0} = p_i$, $q_{ij} = q_{i,j-1} \times q_{i+2^{j-1},j-1}$ 。

首先计算 q_{ij} , 当 u 除以每个 q_{ij} 时, 可以得到余数 u_{ij} 。期望的结果是 u_0 。详细的程序见图 8-3。□

```

begin
1.   for i ← 0 until k-1 do  $q_{i0} \leftarrow p_i$ ;
2.   for j ← 1 until t-1 do
3.     for i ← 0 step  $2^j$  until k-1 do
4.        $q_{ij} \leftarrow q_{i,j-1} * q_{i+2^{j-1},j-1}$ ;
5.        $u_{i0} \leftarrow u$ ;
6.       for j ← t step -1 until 1 do
7.         for i ← 0 step  $2^j$  until k-1 do
8.           begin
9.              $u_{i,j-1} \leftarrow \text{REMAINDER}(u_{ij}/q_{i,j-1})$ ;
10.             $u_{i+2^{j-1},j-1} \leftarrow \text{REMAINDER}(u_{ij}/q_{i+2^{j-1},j-1})$ ;
11.          end;
12.        for i ← 0 until k-1 do  $u_i \leftarrow u_{i0}$ 
13.      end
14.    end
15.  end
end

```

图 8-3 余数计算

定理 8.8 算法 8.4 正确计算 u_0 。

证明: 这个定理的证明和定理 7.3 的证明是平行的, 在定理 7.3 中多项式在 n 次单位根下计值。通过对 j 归纳很容易知道第 4 行对 q_{ij} 的计算是正确的。然后, 通过对第 6 行 j 的逆向归纳可以证明 $u_{ij} = u \bmod q_{ij}$ 。第 8 和 9 行使这个过程变得很容易。假设 $j=0$, 那么有 $u_i = u \bmod p_i$ 。证明的详细过程留作练习。□

定理 8.9 如果每个 p_i 最多需要 b 位来表示, 那么算法 8.4 要 $O_b(M(bk) \log k)$ 时间。

证明: 很容易可以看到第 3~4 行以及第 7~9 行的循环需要花费的时间最多, 每个需要 $O_b(2^{t-j}M(2^{j-1}b))$ 时间^①。假设 $M(an) \geq aM(n)$, 其中 $a \geq 1$, 这些循环花费的代价在 $O_b(M(2^t b)) = O_b(M(kb))$ 范围内。因为每个循环最多循环 $t = \log k$ 次, 由此可以得到结果。□

推论: 如果每个模 p_0, p_1, \dots, p_{k-1} 需要 b 位表示, 那么余数的计算最多可以在 $O_b(bk \log k \log bk \log \log bk)$ 时间内完成。

8.5 多项式模算术和多项式计值

多项式情况下的结果类似于整数。假设 p_0, \dots, p_{k-1} 是多项式, $p = \prod_{i=0}^{k-1} p_i$, 那么每个多项式 u 可以用一系列余数 u_0, u_1, \dots, u_{k-1} 表示, 其中 u_i 是 u 除以 p_i 得到的余数。也就是说, u_i 是满足 $\text{DEG}(u_i) < \text{DEG}(p_i)$ 和 $u = p_i q_i + u_i$ 的唯一多项式。这里, 完全模拟整数模算术, 把求余数记为 $u_i = u \bmod p_i$ 。

类似于引理 8.1, 我们将证明, 如果 p_i 两两互质, u 的次数小于 p 的次数, 例如 $\text{SIZE}(u) < \text{SIZE}(p)$, 那么 $u \leftrightarrow (u_0, u_1, \dots, u_{k-1})$ 是一一对应的。更重要的是, 计算余数的算法 8.4 在 p_i 是多项式而不是整数的时候仍然适用。我们不考虑 $b(p_i)$ 的位数, 而考虑多项式 p_i 的最高次数。当然, 复杂度也是根据算术操作步数, 而不是位操作数衡量。通过这些改变, 可以得到类比于定理 8.9 的结果。

定理 8.10 把算法 8.4 转换到多项式域, 计算出多项式 u 关于多项式 p_0, p_1, \dots, p_{k-1} 的余数可以在 $O_A(M(dk) \log k)$ 时间内完成, 其中 d 是 p_i 次数的上界, 并且 u 的次数少于 $\prod_{i=0}^{k-1} p_i$ 。

① 回忆一下: $D(n)$ 和 $M(n)$ 在实质上是相同的。

证明: 与定理 8.9 类似, 留作练习。□

推论 1 计算多项式 u 关于多项式 p_0, p_1, \dots, p_{k-1} 的余数, 需要的时间最多为 $O_A(dk \log k \log dk)$, 其中 d 是 p_i 次数的上界, 并且 u 的次数少于 $\prod_{i=0}^{k-1} p_i$ 。

例 8.5 考虑下面 4 个多项式的模

$$p_0 = x - 3$$

$$p_1 = x^2 + x + 1$$

$$p_2 = x^2 - 4$$

$$p_3 = 2x + 2$$

假设 $u = x^5 + x^4 + x^3 + x^2 + x + 1$ 。首先计算乘积

$$p_0 p_1 = x^3 - 2x^2 - 2x - 3,$$

$$p_2 p_3 = 2x^3 + 2x^2 - 8x - 8$$

然后计算

$$u' = u \bmod p_0 p_1 = 28x^2 + 28x + 28$$

$$u'' = u \bmod p_2 p_3 = 21x + 21$$

也就是 $u = p_0 p_1 (x^2 + 3x + 9) + 28x^2 + 28x + 28$ 和 $u = p_2 p_3 \left(\frac{1}{2}x^2 + \frac{5}{2} \right) + 21x + 21$ 。

下一步计算

$$u \bmod p_0 = u' \bmod p_0 = 364$$

$$u \bmod p_1 = u' \bmod p_1 = 0$$

$$u \bmod p_2 = u'' \bmod p_2 = 21x + 21$$

$$u \bmod p_3 = u'' \bmod p_3 = 0 \quad \square$$

注意, 7.2 节中的 FFT 算法实际上是这个算法的具体实现, 其中多项式 p_0, p_1, \dots, p_{k-1} 是 $x - \omega^0, x - \omega^1, \dots, x - \omega^{n-1}$ 。FFT 算法利用 $p_i = x - \omega^i$ 这一事实。因为 p_i 的顺序, 每个乘积都有 x 减去 ω 的幂这种形式, 所以除法操作非常容易。

和我们在 7.2 节中看到的一样, 如果 p_i 是一次多项式 $x - a$, 那么 u 模 p_i 为 $u(a)$ 。因此, 所有 p_i 的次数都是 1 的情况非常重要。对定理 8.10 有如下推论。

推论 2 n 次多项式可以在 $O_A(n \log^2 n)$ 时间内在 n 个点进行计值。

证明: 为了在 a_0, a_1, \dots, a_{n-1} 这 n 个点对 $u(x)$ 计值, 当 $0 \leq i < n$ 时, 计算 $u(x)$ 模 $(x - a_i)$ 。因为 d 为 1, 由推论 1 得知这个计值过程需要 $O_A(n \log^2 n)$ 时间。

8.6 中国余数

现在讨论如何把整数从模表示转化为基数表示^①。假设有互质模 p_0, p_1, \dots, p_{k-1} 和余数 u_0, u_1, \dots, u_{k-1} , 其中 $k = 2^t$, 希望找到一个整数 u , 使 $u \leftrightarrow (u_0, u_1, \dots, u_{k-1})$ 。对于整数, 具体操作将类似于多项式的拉格朗日插值公式。

引理 8.2 假设 c_i 是所有 p_j (除了 p_i) 的乘积 (即 $c_i = p/p_i$, 其中 $p = \prod_{j=0}^{k-1} p_j$)。

设 d_i 是 c_i^{-1} 模 p_i (也就是说 $d_i c_i \equiv 1 \bmod p_i$, 且 $0 \leq d_i < p_i$)。那么

$$u \equiv \sum_{i=0}^{k-1} c_i d_i u_i \bmod p \quad (8-20)$$

证明: 因为 p_j 是互质的, 所以 d_i 存在并且唯一 (习题 7.9)。又因为对于 $j \neq i$, c_i 可以被 p_j 整

① 这就是著名的中国余数问题, 因为这个过程中国人在 2000 年前就发现了。

除。所以, 如果 $j \neq i$, 则 $c_i d_i u_i \equiv 0 \pmod{p_j}$ 。所以

$$\sum_{i=0}^{k-1} c_i d_i u_i \equiv c_j d_j u_j \pmod{p_j}$$

因为 $c_j d_j \equiv 1 \pmod{p_j}$, 所以

$$\sum_{i=0}^{k-1} c_i d_i u_i \equiv u_j \pmod{p_j}$$

又因为 p_j 除 p , 所以, 即使所有的算术都模 p , 关系仍然成立。这样式(8-20)成立。□

问题是如何有效地计算式(8-20)。开始时, 使用试错法如何从 p_i 计算 d_i 是很难搞清楚的。但稍后将会看到这个计算任务并不难, 8.8 节将给出解决该问题的欧几里得算法和 8.10 节将给出算法的快速实现。然而在本节, 我们仅研究中国余数问题的“预备好”的形式 (preconditioned form)。如果对于一定数量的问题, 它的一部分输入是固定的, 那么所有依赖于固定输入部分的常量可以预先计算, 并且作为输入的一部分。如果所有这些常量的预计算完成, 那么该算法称为预备好的。

对于预备好的中国余数算法, 输入不但是模和 p_i , 还包括逆 d_i 。这种想法并非不现实的。如果经常使用中国余数算法转化由固定模集表示的数, 那么预先计算所有在算法中用到的模的函数是很合理的。在定理 8.21 的推论中将看到, 只要涉及到数量级, 对于算法是否预备好只产生稍微的不同。

看式(8-20), 可以注意到项 $c_i d_i u_i$ 随着 i 的改变有许多公共的因子。例如, 当 $i \geq k/2$ 时, $c_i d_i u_i$ 有 $p_0 p_1 \cdots p_{k/2-1}$ 作为因子; 当 $i < k/2$ 时, 有 $p_{k/2} p_{k/2+1} \cdots p_{k-1}$ 作为因子。所以, 可以把式(8-20)写成

$$u = \left(\sum_{i=0}^{k/2-1} c_i' d_i u_i \right) \times \prod_{i=k/2+1}^{k-1} p_i + \left(\sum_{i=k/2+1}^{k-1} c_i'' d_i u_i \right) \times \prod_{i=0}^{k/2-1} p_i$$

其中 c_i' 是 $p_0, p_1, \dots, p_{k/2-1}$ 的乘积 (其中不包括 p_i), c_i'' 是 $p_{k/2}, p_{k/2+1}, \dots, p_{k-1}$ 的乘积 (其中不包括 p_i)。这实际上建议使用类似于余数计算中的分治法。计算下面的乘积

$$q_{ij} = \prod_{m=i}^{i+2^j-1} p_m$$

(和 8.4 中类似) 然后计算整数

$$s_{ij} = \sum_{m=i}^{i+2^j-1} q_{ij} d_m u_m / p_m$$

如果 $j=0$, 那么 $s_{i0} = d_i u_i$ 。如果 $j>0$, 则通过下面的公式计算 s_{ij}

$$s_{ij} = s_{i,j-1} q_{i+2^{j-1}, j-1} + s_{i+2^{j-1}, j-1} q_{i, j-1}$$

最后, 计算出 $s_{0k} = u$, 也就是式(8-20)。

算法 8.5 快速中国余数预备好算法。

输入:

1. 互质整数模 p_0, p_1, \dots, p_{k-1} , 其中存在 t , 使 $k=2^t$ 。
2. “逆”的集合 d_0, d_1, \dots, d_{k-1} , 满足 $d_i = (p/p_i)^{-1} \pmod{p_i}$, 其中 $p = \prod_{i=0}^{t-1} p_i$ 。
3. 余数序列 $(u_0, u_1, \dots, u_{k-1})$ 。

输出: 唯一的整数 u , $0 \leq u < p$, 满足 $u \leftrightarrow (u_0, u_1, \dots, u_{k-1})$ 。

方法: 首先如算法 8.4[○], 计算 $q_{ij} = \prod_{m=i}^{i+2^j-1} p_m$ 。然后执行图 8-4 中的程序, 其中 s_{ij} 是

$$\sum_{m=i}^{i+2^j-1} q_{ij} d_m u_m / p_m$$

□

○ 注意, q_{ij} 仅是 p_i 的函数。因为允许预先准备, 所以把它们包含在输入里, 而不是计算得到。然而, 很容易证明运行时间的数量级不受 q_{ij} 是否预计算影响。

```

begin
1.   for i ← 0 until k-1 do  $s_{i0} \leftarrow d_i * u_i$ ;
2.   for j ← 1 until t do
3.       for i ← 0 step  $2^j$  until k-1 do
4.            $s_{ij} \leftarrow s_{i,j-1} * q_{i+2^{j-1},j-1} + s_{i+2^{j-1},j-1} * q_{i,j-1}$ ;
5.   write  $s_{kt}$  modulo  $q_{kt}$ 
end

```

图 8-4 由模表示计算整数的程序

例 8.6 假设 p_0, p_1, p_2, p_3 为 2, 3, 5, 7, (u_0, u_1, u_2, u_3) 为 (1, 2, 4, 3)。对于 $0 \leq i < 4$ 有 $q_0 = p_i$, $q_{01} = 6$, $q_{21} = 35$, $q_{02} = p = 210$ 。记

因为 $1 * 105 \equiv 1 \pmod{2}$, 所以 $d_0 = (3 * 5 * 7)^{-1} \pmod{2} = 1$

因为 $1 * 70 \equiv 1 \pmod{3}$, 所以 $d_1 = (2 * 5 * 7)^{-1} \pmod{3} = 1$

因为 $3 * 42 \equiv 1 \pmod{5}$, 所以 $d_2 = (2 * 3 * 7)^{-1} \pmod{5} = 3$

因为 $4 * 30 \equiv 1 \pmod{7}$, 所以 $d_3 = (2 * 3 * 5)^{-1} \pmod{7} = 4$

如此, 第一行的效果就是计算

$$s_{00} = 1 * 1 = 1, s_{10} = 1 * 2 = 2$$

$$s_{20} = 3 * 4 = 12, s_{30} = 4 * 3 = 12$$

然后执行第 3~4 行的循环, 且 $j=1$ 。其中 i 在 0 和 2 取值, 所以计算

$$s_{01} = s_{00} * q_{10} + s_{10} * q_{00} = 1 * 3 + 2 * 2 = 7,$$

$$s_{21} = s_{20} * q_{30} + s_{30} * q_{20} = 12 * 7 + 12 * 5 = 144.$$

下一步在 $j=2$ 时执行第 3~4 行的循环, i 的取值仅为 0。计算

$$s_{02} = s_{01} * q_{21} + s_{21} * q_{01} = 7 * 35 + 144 * 6 = 1109.$$

第 5 行的结果是 1109 模 210, 为 59。可以验证, 59 模 2, 3, 5, 7 得到的余数分别为 1, 2, 4, 3。图 8-5 用图说明整个计算过程。

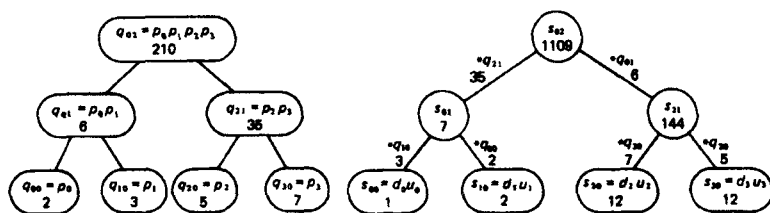


图 8-5 例 8.6 的计算过程

定理 8.11 算法 8.5 正确地计算 u , 其满足 $u \leftrightarrow (u_0, u_1, \dots, u_{k-1})$ 。

证明: 对 j 进行基本归纳, 证明 s_{ij} 取得了预期值, 也就是

$$s_{ij} = \sum_{m=i}^{i+2^j-1} q_{jm} d_m u_m / p_m$$

由引理 8.2 马上可以得到算法的正确性, 也就是说从公式 (8-20) 的正确性马上可以得证。

定理 8.12 假设给定 k 个互质的整数模 p_0, p_1, \dots, p_{k-1} 和余数序列 $(u_0, u_1, \dots, u_{k-1})$ 。如果每个 p_i 要求至多 b 位, 那么存在一个预备好算法能够在 $O_b(M(bk) \log k)$ 时间内计算 u , 使 $0 \leq u < p = \prod_{i=0}^{k-1} p_i$ 和 $u \leftrightarrow (u_0, u_1, \dots, u_{k-1})$, 其中 $M(n)$ 表示两个 n 位整数相乘需要的时间。

证明: q_j 的计算需要 $O_b(M(bk)\log k)$ 时间^①。对算法主体进行分析, 因为 s_j 是 2^j 个项的和, 且每项都是 $2^j + 1$ 个小于等于 b 位的整数之和, 所以 s_j 需要至多 $b2^j + b + j$ 位。每项要求不多于 $b(2^j + 1)$ 位, 2^j 个这样的项要求最多 $b(2^j + 1) + \log(2^j) = b2^j + b + j$ 位。那么第 4 行需要的时间为 $O_b(M(b2^j))$ 。第 3~4 行的循环对于一个固定的 j 需要执行 $k/2^j$ 次, 所以整个循环需要的时间为

$$O_b\left(\frac{k}{2^j}M(b2^j)\right)$$

由关于 $M(n)$ 的增长通常假设, 整个循环需要的时间上界为 $O_b(M(bk))$ 。因为第 2~4 行的循环重复了 $\log k$ 次, 所以总的代价为 $O_b(M(bk)\log k)$ 时间。第 5 行所需要的时间花费很显然少于这个值。□

推论: 在 k 个 b 位整数模的中国余数预备好算法至多需要 $O_b(bk \log k \log bk \log \log bk)$ 步计算。

8.7 中国余数和多项式的插值

很明显, 前一节中有关多项式模 p_0, p_1, \dots, p_{k-1} 的结果对于整数情况仍然成立。所以, 我们可以得到下面的定理和推论。

定理 8.13 假设 $p_0(x), p_1(x), \dots, p_{k-1}(x)$ 是次数最多为 d 的多项式, 而 $M(n)$ 是计算两个 n 次多项式需要的算术操作步数。对于给定多项式 $u_0(x), u_1(x), \dots, u_{k-1}(x)$, 其中对于所有的 $0 \leq i < k$, $u_i(x)$ 的次数少于 $p_i(x)$ 的次数, 存在一个计算唯一多项式 $u(x)$ 的预先准备算法, 其能够在 $O_A(M(dk)\log k)$ 时间内完成计算, 其中, $u(x)$ 的次数小于 $p(x) = \prod_{i=0}^{k-1} p_i(x)$ 的次数, 且 $u(x) \leftrightarrow (u_0(x), u_1(x), \dots, u_{k-1}(x))$ 成立。

证明: 与算法 8.5 和定理 8.12 类似。□

推论: 对于多项式中国余数问题, 存在一个 $O_A(dk \log k \log dk)$ 时间的预先准备算法。

当所有模的次数都是 1 时, 可以得到一个非常重要的特殊例子。如果 $p_i = x - a_i$, 其中 $0 \leq i < k$, 那么 u_i 的余数是常数, 也就是次数为 0 的多项式。如果 $u(x) \equiv u_i \pmod{x - a_i}$, 那么 $u(x) = q(x)(x - a_i) + u_i$, 所以有 $u(a_i) = u_i$ 。那么次数小于 k 且满足 $u(x) \leftrightarrow (u_0, u_1, \dots, u_{k-1})$ 的唯一多项式 $u(x)$ 是唯一的次数小于 k 且对于所有 $0 \leq i < k$ 满足 $u(a_i) = u_i$ 的多项式。另外一个解释是, $u(x)$ 是通过在 (a_i, u_i) 点插值得到的多项式, 其中 $0 \leq i < k$ 。

因为插值是一种非常重要的多项式操作, 即使不允许预先准备, 对 k 点的插值操作也可以在 $O_A(k \log^2 k)$ 时间内完成是很值得高兴的事。就像下一个引理中阐述, 这是由于在这种特殊情况下, 可以很快地计算出公式(8-20)中的系数 d_i ^②。

引理 8.3 假设 $p_i(x) = x - a_i$, $0 \leq i < k$, 其中 a_i 是明显的(也就是说 $p_i(x)$ 是两两互质的)。设 $p(x) = \prod_{j=0}^{k-1} p_j(x)$, $c_i(x) = p(x)/p_i(x)$, $d_i(x)$ 是常量多项式, 满足 $d_i(x)c_i(x) \equiv 1 \pmod{p_i(x)}$ 。那么 $d_i(x) = 1/b$, 其中 $b = \frac{d}{dx}p(x) \big|_{x=a_i}$ 。

证明: 可以写 $p(x) = c_i(x)p_i(x)$, 所以

$$\frac{d}{dx}p(x) = p_i(x) \frac{d}{dx}c_i(x) + c_i(x) \frac{d}{dx}p_i(x) \quad (8-21)$$

现在 $dp_i(x)/dx = 1$ 和 $p_i(a_i) = 0$, 所以

$$\frac{dp(x)}{dx} \big|_{x=a_i} = c_i(a_i) \quad (8-22)$$

① 因为 $D(n)$ 和 $M(n)$ 本质上是相同的函数, 选择使用 $M(n)$ 。

② 和 8.6 节提到的一样, 在通常情况下, 这个任务事实上并不难。然而一般情况下, 需要下一节中的机制。

注意, $d_i(x)$ 满足 $d_i(x)c_i(x) \equiv 1 \pmod{(x-a_i)}$ 的性质, 所以, 存在 $q_i(x)$ 使 $d_i(x)c_i(x) = q_i(x)(x-a_i) + 1$ 成立。现在 $d_i(a_i) = 1/c_i(a_i)$ 。由于 $d_i(x)$ 是个常量, 所以由式(8-22)马上得证此引理。□

定理 8.14 在没有预先准备的情况下, 可以在 $O_A(k \log^2 k)$ 时间内通过 k 点插值一个多项式。

证明: 由引理 8.3, d_i 的计算等同于对 $(k-1)$ 次多项式在 k 个点上求导数。先计算 $p_0 p_1, p_2 p_3, \dots$ 的乘积, 然后计算 $p_0 p_1 p_2 p_3, p_4 p_5 p_6 p_7, \dots$ 等等的乘积, 可以在 $O_A(k \log^2 k)$ 时间内得到多项式 $p(x) = \prod_{i=0}^{k-1} p_i(x)$ 。而 $p(x)$ 的导数可以在 $O_A(k)$ 步内计算出。由定理 8.10 的推论 2 可知, 导数的计算要求 $O_A(k \log^2 k)$ 时间。当 $d=1$ 时, 从推论到定理 8.13 可以得到定理。□

例 8.7 在 $(1, 2), (2, 7), (3, 4)$ 和 $(4, 8)$ 这些点对一个多项式进行插值操作。当 $0 \leq i < 4$ 时, $a_i = i+1, u_0 = 2, u_1 = 7, u_2 = 4$ 和 $u_3 = 8$ 。所以 $p_i(x) = x - i - 1, p(x) = \prod_{i=0}^3 p_i(x)$ 等于 $x^4 - 10x^3 + 35x^2 - 50x + 24$ 。下一步, $dp(x)/dx = 4x^3 - 30x^2 + 70x - 50$, 在 $1, 2, 3, 4$ 上的值分别为 $-6, +2, -2, +6$ 。由此, d_0, d_1, d_2 和 d_3 的值分别为 $-\frac{1}{6}, +\frac{1}{2}, -\frac{1}{2}$ 和 $+\frac{1}{6}$ 。利用快速中国余数算法(算法 8.5)对多项式重新计算, 可以得到

$$s_{01} = d_0 u_0 p_1 + d_1 u_1 p_0 = \left(-\frac{1}{6}\right)(2)(x-2) + \left(\frac{1}{2}\right)(7)(x-1) = \frac{19}{6}x - \frac{17}{6}$$

$$s_{21} = d_2 u_2 p_3 + d_3 u_3 p_2 = \left(-\frac{1}{2}\right)(4)(x-4) + \left(\frac{1}{6}\right)(8)(x-3) = -\frac{2}{3}x + 4$$

那么

$$s_{02} = u(x) = s_{01} q_{21} + s_{21} q_{01} = \left(\frac{19}{6}x - \frac{17}{6}\right)(x^2 - 7x + 12) + \left(-\frac{2}{3}x + 4\right)(x^2 - 3x + 2)$$

$$u(x) = \frac{5}{2}x^3 - 19x^2 + \frac{89}{2}x - 26$$

与第 7 章所提到的一样, 我们可以通过对多项式的 n 点计值, 计算出这些点上的算术值, 再利用这些值对多项式进行插值操作来进行多项式算术计算, 例如加、减、乘。如果得到的结果是 $n-1$ 次或更低次多项式, 那么这个技术得到的结果是正确的。

FFT 就是按照这种做法, 其中所选择的点是 $\omega^0, \omega^1, \dots, \omega^{n-1}$ 。在这种情况下, 由于 ω 的幂属性和为这些幂选定的顺序, 使得计值和插值算法非常简单。然而, 值得注意的是, 可以用任意点集来替代 ω 的幂。需要一次“变换”, 这需要 $O_A(n \log^2 n)$ 时间, 而不是 $O_A(n \log n)$ 时间, 计算和反向。

8.8 最大公因子和欧几里得算法

定义 假设 a_0 和 a_1 是正整数, 如果满足下面两个条件, 则正整数 g 称为 a_0 和 a_1 的最大公因子, 记为 $\text{GCD}(a_0, a_1)$,

1. g 可以整除 a_0 和 a_1 ,
2. 所有 a_0 和 a_1 的因子都可以整除 g 。

很容易得到, 如果 a_0 和 a_1 是正整数, 那么 g 是唯一的。例如, $\text{GCD}(57, 33)$ 等于 3。

用来计算 $\text{GCD}(a_0, a_1)$ 的欧几里得算法是计算余数序列 a_0, a_1, \dots, a_k , 其中对于所有的 $i \geq 2$, a_i 是 a_{i-1} 除以 a_{i-2} 的非零余数, 且 a_{k-1} 除以 a_k 的余数为 0 (例如 $a_{k+1} = 0$)。那么 $\text{GCD}(a_0, a_1) = a_k$ 。

例 8.8 在上面的例子中, $a_0 = 57, a_1 = 33, a_2 = 24, a_3 = 9, a_4 = 6$ 和 $a_5 = 3$, 由此, $k=5$ 且 $\text{GCD}(57, 33) = 3$ 。□

定理 8.15 欧几里得算法正确计算 $\text{GCD}(a_0, a_1)$ 。

证明: 算法计算 $a_{i+1} = a_{i-1} - q_i a_i$, 其中 $1 \leq i < k, q_i = \lfloor a_{i-1}/a_i \rfloor$ 。因为 $a_{i+1} < a_i$, 所以算法肯

定能够终止。而且任何 a_{i-1} 和 a_i 的因子都是 a_{i+1} 的因子, 任何 a_i 和 a_{i+1} 的因子也是 a_{i-1} 的因子。所以 $\text{GCD}(a_0, a_1) = \text{GCD}(a_1, a_2) = \cdots = \text{GCD}(a_{k-1}, a_k)$ 。因为 $\text{GCD}(a_{k-1}, a_k)$ 是 a_k , 所以结论成立。□

欧几里得算法的扩展不但可以用来找 a_0 和 a_1 的最大公因子, 还可以计算整数 x 和 y , 使它们满足 $a_0x + a_1y = \text{GCD}(a_0, a_1)$ 。算法如下。

算法 8.6 扩展的欧几里得算法。

输入: 正整数 a_0 和 a_1 。

输出: $\text{GCD}(a_0, a_1)$ 和满足 $a_0x + a_1y = \text{GCD}(a_0, a_1)$ 的整数 x 和 y 。

方法: 执行图 8-6 中的程序。□

```

begin
1.   $x_0 \leftarrow 1; y_0 \leftarrow 0; x_1 \leftarrow 0; y_1 \leftarrow 1; i \leftarrow 1;$ 
2.  while  $a_i$  除以  $a_{i-1}$  余数不为 0 do
      begin
3.       $q \leftarrow \lfloor a_{i-1}/a_i \rfloor;$ 
4.       $a_{i+1} \leftarrow a_{i-1} - q * a_i;$ 
5.       $x_{i+1} \leftarrow x_{i-1} - q * x_i;$ 
6.       $y_{i+1} \leftarrow y_{i-1} - q * y_i;$ 
7.       $i \leftarrow i + 1$ 
      end
8.  write  $a_i$ ; write  $x_i$ ; write  $y_i$ 
end

```

图 8-6 扩展的欧几里得算法

例 8.9 如果 $a_0 = 57$, $a_1 = 33$, 则可以得到下面的 a_i , x_i 和 y_i 值。

i	a_i	x_i	y_i
0	57	1	0
1	33	0	1
2	24	1	-1
3	9	-1	2
4	6	3	-5
5	3	-4	7

注意 $57 \times (-4) + 33 \times 7 = 3$ 。□

由于 a_i 很清楚地形成了余数序列, 所以很显然, 算法 8.6 正确地计算出了 $\text{GCD}(a_0, a_1)$ 。算法 8.6 中计算的 x_i 和 y_i 有一个非常重要的属性将在下面这个引理中讲到。

引理 8.4 在算法 8.6 中, 对于 $i \geq 0$, 有

$$a_0x_i + a_1y_i = a_i \quad (8-23)$$

证明: 当 $i=0$ 和 $i=1$ 时, 由算法 8.6 的第 1 行可知, 公式(8-23)成立。假设式(8-23)对 $i-1$ 和 i 成立, 那么由算法第 5 行有 $x_{i+1} = x_{i-1} - qx_i$, 由算法第 6 行有 $y_{i+1} = y_{i-1} - qy_i$ 。所以

$$a_0x_{i+1} + a_1y_{i+1} = a_0x_{i-1} + a_1y_{i-1} - q(a_0x_i + a_1y_i) \quad (8-24)$$

由归纳假设和式(8-24), 有

$$a_0x_{i+1} + a_1y_{i+1} = a_{i-1} - qa_i$$

因为由第 4 行可得 $a_{i-1} - qa_i = a_{i+1}$, 所以结论得证。□

注意, 虽然算法第 3 行保证算法 8.6 确实计算出了 $\text{GCD}(a_0, a_1)$, 但是引理 8.4 并不依赖于算法第 3 行中 q 的计算方法。把这些观察结果综合起来可得到下面的定理。

定理 8.16 算法 8.6 计算 $\text{GCD}(a_0, a_1)$, 数 x 和 y 满足 $a_0x + a_1y = \text{GCD}(a_0, a_1)$ 。

证明: 应用引理 8.4 的基本练习。 \square

现在介绍一些在下一节中用到的概念。

定义: 假设 a_0 和 a_1 是整数, 其余数序列为 a_0, a_1, \dots, a_k 。设 $q_i = \lfloor a_{i-1}/a_i \rfloor$, 其中 $1 \leq i \leq k$ 。

定义 2×2 矩阵 $R_{ij}^{(a_0, a_1)}$ 或 R_{ij} , 其中 $0 \leq i \leq j \leq k$, a_0 和 a_1 事先确定如下:

$$1. \text{ 对于 } i \geq 0, R_{ii} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix};$$

$$2. \text{ 如果 } j > i, \text{ 那么 } R_{ij} = \begin{bmatrix} 0 & 1 \\ 1 & -q_j \end{bmatrix} * \begin{bmatrix} 0 & 1 \\ 1 & -q_{j-1} \end{bmatrix} * \dots * \begin{bmatrix} 0 & 1 \\ 1 & -q_{i+1} \end{bmatrix}.$$

例 8.10 假设 $a_0 = 57, a_1 = 33$, 其余数序列为 57, 33, 24, 9, 6, 3, 商为 $q_i (1 \leq i \leq 4)$ 分别为 1, 1, 2, 1。那么

$$R_{04}^{(57, 33)} = \begin{bmatrix} 0 & 1 \\ 1 & -1 \end{bmatrix} * \begin{bmatrix} 0 & 1 \\ 1 & -2 \end{bmatrix} * \begin{bmatrix} 0 & 1 \\ 1 & -1 \end{bmatrix} * \begin{bmatrix} 0 & 1 \\ 1 & -1 \end{bmatrix} = \begin{bmatrix} 3 & -5 \\ -4 & 7 \end{bmatrix} \quad \square$$

下面引理给出这些矩阵的两个属性。

引理 8.5

$$a) \begin{bmatrix} a_j \\ a_{j+1} \end{bmatrix} = R_{ij} \begin{bmatrix} a_i \\ a_{i+1} \end{bmatrix}, \text{ 其中 } i < j < k,$$

$$b) R_{0j} = \begin{bmatrix} x_j & y_j \\ x_{j+1} & y_{j+1} \end{bmatrix}, \text{ 其中 } 0 \leq j < k,$$

其中 a_i, x_i 和 y_i 和扩展欧几里得算法中的定义相同。

证明: 基本归纳练习。 \square

应该注意到, 经过下面修改之后, 本节中的所有结论不仅对整数, 对一元多项式也成立。很显然, 两个整数的最大公因子是唯一的; 而对多项式而言, 其最大公因子只有忽略常量乘数才是唯一的。也就是说, 如果 $g(x)$ 整除多项式 $a_0(x)$ 和 $a_1(x)$, 这两个多项式的其他因子也能够整除 $g(x)$, 那么对于所有的常量 $c \neq 0$, $cg(x)$ 都满足这一性质。我们将满足于找到能够整除 $a_0(x)$ 和 $a_1(x)$ 并且能够被所有其他因子(多项式)整除的任意多项式即可。为了唯一性, 可以(但不)坚持最大的公因子是首项系数为 1 的, 即最高次项的系数为 1 的多项式。

8.9 多项式 GCD 的渐近快速算法

针对整数的算法需要处理一些特殊的细节, 所以, 我们首先讨论多项式的最大公因子算法。假设 $a_0(x)$ 和 $a_1(x)$ 是两个需要计算最大公因子的多项式, 且假设 $\text{DEG}(a_1(x)) < \text{DEG}(a_0(x))$ 。这个条件可以通过如下方式强制满足。如果 $\text{DEG}(a_0) = \text{DEG}(a_1)$, 那么把多项式 a_0 和 a_1 用 a_1 和 a_0 模 a_1 替换, 也就是余数序列的第 2 和第 3 项, 再从这里开始处理。

把问题分为两部分。第一部分设计算法获取余数序列的最后一项, 该项的次数大于多项式 a_0 次数的一半。形式化地说, 假设 $l(i)$ 是使得 $\text{DEG}(a_{l(i)}) > i$ 和 $\text{DEG}(a_{l(i)+1}) \leq i$ 成立的唯一整数。注意, 如果 a_0 的次数为 n , 那么在假设 $\text{DEG}(a_1) < \text{DEG}(a_0)$ 的情况下有 $l(i) \leq n - i - 1$, 原因在于, 对于所有的 $i \geq 1$ 有 $\text{DEG}(a_i) \leq \text{DEG}(a_{i-1}) - 1$ 成立。

现在引入一个递归过程 HGCD(half GCD), 其输入 a_0 和 a_1 , 产生矩阵 R_{0j} (见 8.8 节), 其中 $n = \text{DEG}(a_0) > \text{DEG}(a_1)$, $j = l(n/2)$, 也就是说, a_j 这个余数序列最后一项的次数大于 a_0 次数的一半。

HGCD 算法背后的规则是次数为 d_1 和 d_2 的多项式的商仅依赖于被除数的头 $2(d_1 - d_2) + 1$ 项和除数的头 $d_1 - d_2 + 1$ 项, 其中 $d_1 > d_2$ 。HGCD 定义如图 8-7 所示。

```

procedure HGCD( $a_0, a_1$ ):
1. if  $\text{DEG}(a_1) \leq \text{DEG}(a_0)/2$  then return  $\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$ 
   else
     begin
2.   令  $a_0 = b_0x^m + c_0$ , 其中  $m = \lfloor \text{DEG}(a_0)/2 \rfloor$  和  $\text{DEG}(c_0) < m$ ;
3.   令  $a_1 = b_1x^m + c_1$ , 其中  $\text{DEG}(c_1) < m$ ;
       comment  $b_0$  和  $b_1$  是  $a_0$  和  $a_1$  的前导项 (leading term)。有
        $\text{DEG}(b_0) = \lfloor \text{DEG}(a_0)/2 \rfloor$  和  $\text{DEG}(b_0) - \text{DEG}(b_1) = \text{DEG}(a_0) - \text{DEG}(a_1)$ ;
4.    $R \leftarrow \text{HGCD}(b_0, b_1)$ ;
5.    $\begin{bmatrix} d \\ e \end{bmatrix} \leftarrow R \begin{bmatrix} a_0 \\ a_1 \end{bmatrix}$ ;
6.    $f \leftarrow d \bmod e$ ;
       comment  $e$  和  $f$  是余数序列中的连续项, 它们的次数最大为  $\lfloor 3m/2 \rfloor$ , 也就是  $a_0$  次数的  $3/4$ ;
7.   令  $e = g_0x^{\lfloor m/2 \rfloor} + h_0$ , 其中  $\text{DEG}(h_0) < \lfloor m/2 \rfloor$ ;
8.   令  $f = g_1x^{\lfloor m/2 \rfloor} + h_1$ , 其中  $\text{DEG}(h_1) < \lfloor m/2 \rfloor$ ;
       comment  $g_0$  和  $g_1$  的次数最多为  $m+1$ ;
9.    $S \leftarrow \text{HGCD}(g_0, g_1)$ ;
10.   $q \leftarrow \lfloor d/e \rfloor$ ;
11.  return  $S * \begin{bmatrix} 0 & 1 \\ 1 & -q \end{bmatrix} * R$ 
     end

```

图 8-7 过程 HGCD

例 8.11 设

$$p_1(x) = x^5 + x^4 + x^3 + x^2 + x + 1$$

和

$$p_2(x) = x^4 - 2x^3 + 3x^2 - x - 7$$

假设要计算 $\text{HGCD}(p_1, p_2)$ 。如果 $a_0 = p_1$, $a_1 = p_2$, 那么在第 2 和第 3 行有 $m=2$, 且

$$b_0 = x^3 + x^2 + x + 1, \quad c_0 = x + 1$$

$$b_1 = x^2 - 2x + 3, \quad c_1 = -x - 7$$

那么, 在第 4 行调用 $\text{HGCD}(b_0, b_1)$ 。在这一步可以检查 R 的值为

$$\begin{bmatrix} 0 & 1 \\ 1 & -(x+3) \end{bmatrix}$$

下一步在第 5 和第 6 行, 计算

$$d = x^4 - 2x^3 + 3x^2 - x - 7$$

$$e = 4x^3 - 7x^2 + 11x + 22$$

$$f = -\frac{3}{16}x^2 - \frac{93}{16}x - \frac{45}{8}.$$

我们发现 $\lfloor m/2 \rfloor = 1$, 所以在第 7 和第 8 行可以得到

$$g_0 = 4x^2 - 7x + 11, \quad h_0 = 22$$

$$g_1 = -\frac{3}{16}x - \frac{93}{16}, \quad h_1 = -\frac{45}{8}$$

这样, 在第 9 行可以得到

$$S = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

在第10行, $\lfloor d(x)/e(x) \rfloor$ 的商 $q(x)$ 为 $\frac{1}{4}x - \frac{1}{16}$ 。所以在第11行, 可以得到如下结果

$$T = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 0 & 1 \\ 1 & -(\frac{1}{4}x - \frac{1}{16}) \end{bmatrix} \begin{bmatrix} 0 & 1 \\ 1 & -(x+3) \end{bmatrix} = \begin{bmatrix} 1 & -(x+3) \\ -(\frac{1}{4}x - \frac{1}{16}) & \frac{1}{4}x^2 + \frac{11}{16}x + \frac{13}{16} \end{bmatrix}$$

注意, 在 p_1 和 p_2 的余数序列中, e 是最后一项次数超过 p_1 一半的多项式, 所以有

$$T \begin{bmatrix} p_1 \\ p_2 \end{bmatrix} = \begin{bmatrix} e \\ f \end{bmatrix}$$

□

现在考虑 HGCD 中第4行中计算的矩阵 R 。 R 可能为

$$\prod_{j=1}^{\lfloor \lceil m/2 \rceil \rfloor} \begin{bmatrix} 0 & 1 \\ 1 & -q_j \end{bmatrix}$$

其中 $q_j(x)$ 是 b_0 和 b_1 的余数序列商的第 j 项。也就是说 $R = R_{0, \lfloor \lceil m/2 \rceil \rfloor}^{(b_0, b_1)}$ 。然而在第5行中, 我们用 R 来表示从 d 和 e 中产生的矩阵 $R_{0, \lfloor \lceil 3m/2 \rceil \rfloor}^{(a_0, a_1)}$, 其中 d 是余数序列中次数大于 $3m/2$ 项中的最后一项。

我们要说明这两种 R 的表示方法都是正确的。也就是说, $R_{0, \lfloor \lceil 3m/2 \rceil \rfloor}^{(a_0, a_1)} = R_{0, \lfloor \lceil m/2 \rceil \rfloor}^{(b_0, b_1)}$ 。类似地, 我们可以说明在第9行中的 S 实现了给它赋值的函数。也就是说,

$$S = R_{0, \lfloor \lceil m/2 \rceil \rfloor}^{(d_0, d_1)} = R_{0, \lfloor \lceil m \rceil \rfloor}^{(e, f)}$$

下一个引理蕴涵了这些结果。

引理 8.6 设

$$f(x) = f_1(x)x^k + f_2(x) \quad (8-25)$$

其中 $\text{DEG}(f_2) < k$, 且设

$$g(x) = g_1(x)x^k + g_2(x) \quad (8-26)$$

其中, $\text{DEG}(g_2) < k$ 。假设 q 和 q_1 分别是 $\lfloor f(x)/g(x) \rfloor$ 和 $\lfloor f_1(x)/g_1(x) \rfloor$ 的商, r 和 r_1 分别是 $f(x) - q(x)g(x)$ 和 $f_1(x) - q_1(x)g_1(x)$ 的余数。如果 $\text{DEG}(f) > \text{DEG}(g)$ 且 $k \leq 2\text{DEG}(g) - \text{DEG}(f)$ [即 $\text{DEG}(g_1) \geq \frac{1}{2}\text{DEG}(f_1)$] , 那么

a) $q(x) = q_1(x)$ 且

b) 对于次数是 $k + \text{DEG}(f) - \text{DEG}(g)$ 或者更高的项, $r(x)$ 和 $r_1(x)x^k$ 相等。

证明: 考虑用普通的除法算法计算 $f(x)$ 除以 $g(x)$, 该算法用 $f(x)$ 的第一项除以 $g(x)$ 的第一项得到商的第一项。商的第一项乘以 $g(x)$, 然后从 $f(x)$ 中被减去, 依此类推。前 $\text{DEG}(g) - k$ 项的产生不依赖于 $g_2(x)$ 。但是, 商只包含次数为 $\text{DEG}(f) - \text{DEG}(g)$ 的项。如果 $\text{DEG}(f) - \text{DEG}(g) \leq \text{DEG}(g) - k$, 也就是 $k \leq 2\text{DEG}(g) - \text{DEG}(f)$, 那么商不依赖于 $g_2(x)$ 。如果 $\text{DEG}(f) - \text{DEG}(g) \leq \text{DEG}(f) - k$, 那么商就不依赖于 $f_2(x)$ 。但是, $\text{DEG}(f) - \text{DEG}(g) \leq \text{DEG}(f) - k$ 是因为 $k \leq 2\text{DEG}(g) - \text{DEG}(f)$ 且 $\text{DEG}(f) > \text{DEG}(g)$ 。所以(a)部分命题成立。对于(b)部分, 类似的推理说明次数为 $\text{DEG}(f) - (\text{DEG}(g) - k)$ 或者更高的项不依赖于 $g_2(x)$ 。与此类似, 次数为 k 或者更高的余数项不依赖于 $f_2(x)$ 。但是, $\text{DEG}(f) - \text{DEG}(g) + k > k$ 。这样, $r(x)$ 和 $r_1(x)x^k$ 中所有项的次数是 $\text{DEG}(f) - \text{DEG}(g) + k$ 或者更高。

引理 8.7 假设 $f(x) = f_1(x)x^k + f_2(x)$ 和 $g(x) = g_1(x)x^k + g_2(x)$, 其中, $\text{DEG}(f_2) < k$ 和 $\text{DEG}(g_2) < k$ 。令 $\text{DEG}(f) = n$ 和 $\text{DEG}(g) < \text{DEG}(f)$ 。那么

$$R_{0, \lfloor \lceil (n+k)/2 \rceil \rfloor}^{(f, g)} = R_{0, \lfloor \lceil (n-k)/2 \rceil \rfloor}^{(f_1, g_1)}$$

也就是说, 至少直到余数的次数不多于 f_1 次数的一半时, (f, g) 和 (f_1, g_1) 余数序列的商都

○ 当然, 这个计算也可以放在第5行后面。

是相同的。

证明：由引理 8.6，商相同，并且对于两个余数序列中对应的余数，足够的高次项也相同。 \square

定理 8.17 设 $a_0(x)$ 和 $a_1(x)$ 是多项式，其中 $\text{DEG}(a_0) = n$ ， $\text{DEG}(a_1) < n$ 。那么 $\text{HGCD}(a_0, a_1) = R_{0, l(n/2)} \circ$

证明：直接对 n 进行归纳，利用定理 8.7 可确保第 4 行中 R 为 $R_{0, l(\lceil 3m/2 \rceil)}^{(a_0, a_1)}$ ，第 9 行中的 S 为 $R_{l(\lceil 3m/2 \rceil) + 1, l(m)}^{(a_0, a_1)}$ \square

定理 8.18 如果参数次数最大为 n ，那么 HGCD 需要 $O_A(M(n) \log n)$ 时间，其中 $M(n)$ 表示两个次数为 n 的多项式相乘需要的时间。

证明：证明 n 是 4 的幂时的结果。因为显然 HGCD 需要的时间是非递减函数，所以定理对所有 n 都适用。如果 $\text{DEG}(a_0)$ 是 4 的幂，那么

$$\text{DEG}(b_0) = \frac{1}{2} \text{DEG}(a_0) \text{ 和 } \text{DEG}(g_0) \leq \frac{1}{2} \text{DEG}(a_0)$$

如此，当 HGCD 输入为 n 次时，则存在常数 c 使时间 $T(n)$ 的上界满足

$$T(n) \leq 2T\left(\frac{n}{2}\right) + cM(n) \quad (8-27)$$

也就是说，HGCD 过程体涉及两次对自身的调用，且调用参数是原参数大小的一半和常数次的其他操作，这些操作的时间复杂度是 $O_A(n)$ 或者是 $O_A(M(n))$ 。式(8-27)的求解大家很熟悉。它的上界是 $c_1 M(n) \log n$ ， c_1 是常数。 \square

下面构造一个完备的求解最大公因子的算法。用 HGCD 计算 $R_{0, n/2}$ ，接着是 $R_{0, 3n/4}$ ，然后为 $R_{0, 7n/8}$ 等等，其中 n 是输入的次数。

算法 8.7 GCD 算法。

输入：多项式 $p_1(x)$ 和 $p_2(x)$ ，其中 $\text{DEG}(p_2) < \text{DEG}(p_1)$ 。

输出：多项式 p_1 和 p_2 的最大公因子 $\text{GCD}(p_1, p_2)$ 。

方法：用调用 $\text{GCD}(p_1, p_2)$ ，其中 GCD 是图 8-8 的递归过程。 \square

```

procedure GCD( $a_0, a_1$ ):
1.  if  $a_1$  整除  $a_0$  then return  $a_1$ 
   else
   begin
2.      $R \leftarrow \text{HGCD}(a_0, a_1)$ ;
3.      $\begin{bmatrix} b_0 \\ b_1 \end{bmatrix} \leftarrow R \circ \begin{bmatrix} a_0 \\ a_1 \end{bmatrix}$ ;
4.     if  $b_1$  整除  $b_0$  then return  $b_1$ 
   else
   begin
5.        $c \leftarrow b_0 \text{ modulo } b_1$ ;
6.       return GCD( $b_1, c$ )
   end
   end
end
    
```

图 8-8 GCD 过程

例 8.12 接着例 8.11。其中 $p_1(x) = x^5 + x^4 + x^3 + x^2 + 1$ 和 $p_2(x) = x^4 - 2x^3 + 3x^2 - x - 7$ 。已经计算得到

$$\text{HGCD}(p_1, p_2) = \begin{bmatrix} 1 & -(x+3) \\ -\left(\frac{1}{4}x - \frac{1}{16}\right) & \frac{1}{4}x^2 + \frac{11}{16}x + \frac{13}{16} \end{bmatrix}$$

那么在第3行计算得到 $b_0 = 4x^3 - 7x^2 + 11x + 22$ 和 $b_1 = -\frac{3}{16}x^2 - \frac{93}{16}x - \frac{45}{8}$ 。我们发现 b_1 不能整除 b_0 。在第5行中,有

$$b_0 \bmod b_1 = 3952x + 3952$$

因为后面一项除以 $-\frac{3}{16}x^2 - \frac{93}{16}x - \frac{45}{8}$, 在第6行中对 GCD 的调用在第1行结束, 并且得到 $3952x + 3952$ 作为结果。当然, $x+1$ 也是 p_1 和 p_2 的最大公因子。□

如果能够证明算法 8.7 是可终止的, 则其正确性是很显然的。算法时间复杂度的分析蕴涵了算法的正确性, 也就是下面这个定理。

定理 8.19 如果 $\text{DEG}(p_i) = n$, 那么算法 8.7 需要 $O_A(M(n) \log n)$ 时间, 其中 $M(n)$ 表示两个 n 次多项式相乘需要的时间。

证明: 不等式

$$T(n) \leq T\left(\frac{n}{2}\right) + c_1 M(n) + c_2 M(n) \log n \quad (8-28)$$

描述了算法 8.7 需要的运行时间, 其中 c_1 和 c_2 是常数。也就是说 b_1 的次数少于 a_0 次数的一半, 所以式(8-28)的第1项表示第6行中递归调用需要的时间。项 $c_1 M(n)$ 表示在第1, 3, 4 和 5 行的除法和乘法操作需要的时间, 最后一项表示在第2行中对 HGCD 的调用需要的时间。很容易就可以得到式(8-28)的解, 它的上界是 $kM(n) \log n$, 其中 k 是常数。□

推论: 两个次数最大为 n 的多项式的 GCD 可以在 $O_A(n \log^2 n)$ 时间内计算得到。

8.10 整数的 GCD

本节中要讨论 GCD 和 HGCD 的修改, 使之适用于整数。为了搞清楚整数情况下的问题, 先看引理 8.6。引理 8.6 说明当计算 n 次和 $n-d$ 次多项式的商时, 不需要次数少于 $n-2d$ 的项。

和引理 8.6 类似, 考虑两个整数 g 和 f , 其中 $f > g$, 把 f 和 g 分别写成 $f = f_1 2^k + f_2$ 和 $g = g_1 2^k + g_2$, 其中 $f_2 < 2^k$ 和 $g_2 < 2^k$ 。替换条件 $\text{DEG}(g_1(x)) \geq \frac{1}{2} \text{DEG}(f_1(x))$, 假定 $f_1 \leq (g_1)^2$ 。那么, 可以设 $f = qg + r$ 和 $f_1 = q_1 g_1 + r_1$ 。把这些公式结合, 得到

$$g_1 2^k (q - q_1) = f_2 + r_1 2^k - q g_2 - r \quad (8-29)$$

因为 $r_1 < g_1$, $f_2 < 2^k$, 并且所有的整数都是非负的, 从式(8-29)中可以很容易地得到 $q - q_1 \leq 0$ 。假设存在 $m \geq 0$, 使 $q = q_1 - m$ 成立。那么从式(8-29)可以得到

$$m g_1 2^k \leq q g_2 + r = (q_1 - m) g_2 + r$$

所以

$$m g = m g_1 2^k + m g_2 \leq q_1 g_2 + r \quad (8-30)$$

因为 $f_1 \leq (g_1)^2$, 有 $q_1 \leq g_1$ 。而且已知 $g_2 < 2^k$ 和 $r < g$, 所以式(8-30)隐含

$$m g < g_1 2^k + g \leq 2g \quad (8-31)$$

式(8-31)马上可以得到当 $m < 2$ 的情况。也就是 $q = q_1$ 或者 $q = q_1 - 1$ 。

在前一种情况中, 不存在问题。另一方面, 如果 $q = q_1 - 1$, 我们不能期望 HGCD 能很好地工作。幸运的是, 只要当商是余数序列的最后一个, 我们可以得到 $q_1 \neq q$, 其中余数序列的矩阵由 HGCD 生成。也就是说, 如果把 $q - q_1 = -1$ 代入式(8-29), 那么有

$$r = f_2 + r_1 2^k - q g_2 + g_1 2^k \quad (8-32)$$

因为 $r < g = g_1 2^k + g_2$, 所以可以从式(8-32)推出 $r_1 2^k + f_2 < g_2(1 + q)$, 或者 $r_1 < 1 + q$, 也就是 $r_1 < q_1$ 。

由此, r_1 是余数序列中跟在 f_1 和 g_1 后的项, 必定比 f_1/g_1 小。因为 $g_1 \geq \sqrt{f_1}$, 所以有 $r_1 < \sqrt{f_1}$, 这意味着 HGCD 将返回一个包含到 f_1/g_1 的商(且没有比这个商更小的项)的矩阵。如果这个矩阵在 HGCD 中的第 5 行中用到, 那么在第 6 行中计算的 f 可能不会比 $a_0^{3/4}$ 少。然而, 因为最后商中有一个唯一的错误, 所以在第 6 行后, 扩展余数序列到有限数量(和 a_0 大小无关)使其能够充分地产生小于 $a_0^{3/4}$ 的序列。在过程中第 5 行后需要对 GCD 进行类似的“修正”。

涉及引理 8.6 的另外一个问题是, 在多项式下, 我们能够说明由考虑多项式统一组成的余数序列和完整的多项式形成的序列的特定的高次项是相等的。即使 $q = q_1$, 类似的结果在整数情况下也成立。然而, 需要把 r 和 $r_1 2^t$ 之间的区别限定在 $2^t(q+1)$ 范围内; 我们不能确定特定位的 r 和 $r_1 2^t$ 都相同。尽管如此, “舍入误差”(rounding error)会引起在 HGCD 第 6 行后和 GCD 第 5 行后需要有限数量的余数序列的补充项。

假设 $M(n)$ 表示 n 位数乘法需要的时间, 扩展余数序列所需要的额外代价限制在 $O_b(M(n))$ 范围内, 所以 HGCD 和 GCD 的时间分析基本不受影响。由此有如下定理。

定理 8.20 如果 $M(n)$ 表示两个 n 位数乘法需要的时间, 那么, 存在一个算法能够在 $O_b(M(n) \log n)$ 时间内找到整数 a_0 和 a_1 的最大公因子 $\text{GCD}(a_0, a_1)$ 。

证明: 基于上面对过程 HGCD 和 GCD 的修改建议, 这个证明过程留作练习。□

推论 能够在 $O_b(n \log^2 n \log \log n)$ 时间内找到整数的最大公因子。

8.11 再论中国余数

现在再来看 GCD 算法如何运用于渐近快速且无需预备好的整数的中国余数算法的设计中。回想前面预备好的算法 8.5, 需要 $O_b(M(bk) \log k)$ 时间才能从 k 个 b 位的模上重构 u 。问题是计算 $d_i = (p/p_i)^{-1}$ 模 p_i , 其中是 p_0, p_1, \dots, p_{k-1} 是模, p 是它们的乘积。

利用分治算法, 可以在 $O_b(M(bk) \log k)$ 时间内通过首先成对计算 p_i 的乘积, 然后计算四个 p_i 的乘积等等, 用这样的方法计算 p 。算法 8.5 的技术使我们能够在无需预先计算的条件下, 在 $O_b(M(bk) \log k)$ 步内计算 $e_i = (p/p_i)$ 模 p_i ($0 \leq i < k$)。剩下的就是确定计算 $d_i = e_i^{-1}$ 模 p_i 需要的时间了。

因为 p/p_i 是除了 p_i 外其他模的乘积, 所以它和 p_i 是互质的。存在一个整数 q , 如果把 p/p_i 表示为 $qp_i + e_i$, 那么 e_i 和 p_i 是互质的, 也就是 $\text{GCD}(e_i, p_i) = 1$ 。由此, 给定 x 和 y , 使 $e_i x + p_i y = 1$, 有 $e_i x \equiv 1$ 模 p_i 。也就有 $x \equiv e_i^{-1} = d_i$ 模 p_i 成立。扩展的欧几里得算法就是计算这样的 x 和 y 。

设计 GCD 过程的目的是用来计算 $\text{GCD}(p_1, p_2)$, 而设计 HGCD 则是计算矩阵 $R_{0, l(n/2)}$ 。由此, 对 GCD 算法进行一些小的修改可以生成 $R_{0, l(n)}$ 。因为它将是矩阵的左上元素, 所以有可能取得 x 。现在可能确定没有预先计算的中国余数算法需要的时间。

定理 8.21 给定 k 个 b 位模, 在整数情况下, 中国余数算法需要的时间是 $O_b(M(bk) \log k) + O_b(kM(b) \log b)$ 。

证明: 通过以上分析, 第 1 项表示 e_i 的计算和执行算法 8.5 需要的时间。因为 x 和 y 的计算都允许 b 位算术的模 p_i 计算, 所以第 2 项表示计算 d_i 的时间。□

推论 无预先准备的中国余数算法需要时间最多为 $O_b(bk \log^2 bk \log \log bk)^\ominus$ 。

8.12 稀疏多项式

前面所使用的一元多项式表示基于多项式 $\sum_{i=0}^{n-1} a_i x^i$ 是稠密的这一假设, 也就是说几乎所有的多项式系数都是非 0 的。在许多应用中假设多项式是稀疏的, 也就是说非 0 系数个数远少于多项式的最高次数。在这种情况下多项式的逻辑表示是数对 (a_i, j_i) 的列表, 由非 0 系数和 x 的对应幂组成。

由于不可能对所有处理稀疏多项式的算术运算技术面面俱到, 这里仅涉及两个比较有趣的理论。第一, 在这种情况下, 使用傅里叶变换做乘法运算是不合理的, 将给出一种合理的计算稀疏多项式乘法的技术。第二, 通过计算 $[p(x)]^4$ 说明稠密和稀疏多项式的算术操作存在着很大的差异。

大多数合理的处理稀疏多项式乘法的已知方法是把多项式 $\sum_{i=1}^n a_i x^{j_i}$ 表示为数对的列表 $(a_1, j_1), (a_2, j_2), \dots, (a_n, j_n)$, 其中假设 j 是不同的且按照降序排列, 即对于 $1 \leq i < n$, $j_i > j_{i+1}$ 。为了把两个这种方法表示的多项式相乘, 先计算数对的乘积并且按照它们的指数(数对中的第二个分量)大小排列这些结果项, 然后把具有相同指数的项合并。不这样做的后果是会增加许多具有相同指数的项。这种情况下, 如果在每一步算术操作都进行合并项操作, 随着越来越多的算术操作, 所需要的代价就会大大超出原来的情况。

如果对排序后的稀疏多项式做乘法, 可以利用两个事实使乘积的排序变得尽可能简单, 其一是它们是有序的, 其二是其中一个的项比另外一个多得多。下面设计一个非形式的算法按照这种方法做稀疏多项式的乘法。

算法 8.8 有序稀疏多项式乘法。

输入: 多项式

$$f(x) = \sum_{i=1}^m a_i x^{j_i} \text{ 和 } g(x) = \sum_{i=1}^n b_i x^{k_i}$$

表示为数对列表

$$(a_1, j_1), (a_2, j_2), \dots, (a_m, j_m) \text{ 和 } (b_1, k_1), (b_2, k_2), \dots, (b_n, k_n),$$

其中 j 和 k 是单调递减的。

输出:

$$\sum_{i=1}^p c_i x^{l_i} = f(x)g(x)$$

用数对列表表示, 其中 l_i 是单调递减的。

方法: 不失一般性, 假设 $m \geq n$ 。

1. 对 $1 \leq i \leq n$, 构造 S_i 序列, 其中第 r 项是 $(a_r b_i, j_r + k_i)$ ($1 \leq r \leq m$)。也就是说 S_i 表示 $f(x)$ 和 $g(x)$ 第 i 项的乘积;

2. 对于 $1 \leq i \leq n/2$, 通过项的合并把 S_{2i-1} 序列和 S_{2i} 序列合并。然后按对把结果序列项合并, 一直重复这一过程直到只剩下一个有序的序列。

定理 8.22 在 $m \geq h$ 时, 算法 8.8 需要 $O(mn \log n)$ 时间^①。

证明: 第 1 步肯定需要的时间为 $O(mn)$ 。第 2 步重复执行 $\lceil \log n \rceil$ 次, 每次需要 $O(mn)$ 时间。

① 注意很有趣的一个现象是 $M(n) = n \log n \log \log n$, 这个数字是从定理 8.21 中可以得到的最好情况, 它和 b, k 无关。

② 注意这里使用 RAM 复杂度而不是算术复杂度, 因为即使 m 和 n 是固定的, 算法 8.8 的程序中分支转移也是固有的。

现在来看算法 8.8 及其时间复杂性如何影响稀疏多项式的算术操作。

例 8.13 现在考虑计算多项式 $p^4(x)$, 其中 $p(x)$ 是有 n 项的多项式, 有稀疏和稠密两种情况。当给定的 $p(x)$ 是稠密多项式时, 那么计算 $p^4(x)$ 的最好方法是通过两次平方的计算。也就是说假设 $M(n)$ 是两个稠密多项式相乘需要的时间为 $cn \log n$, 那么, 可以在 $cn \log n$ 步内计算出 $p^2(x)$, 并且在 $2cn \log 2n$ 步内计算出这一结果的平方, 所以总共需要 $3cn \log n + 2cn$ 步。与之进行比较, 如果按照 $p^4 = p \times (p \times (p \times p))$ 的方式计算 $p^4(x)$, 那么在计算 $p \times p^2$ 需要 $2M(n)$ 步, 计算 $p \times p^3$ 需要 $3M(n)$ 步的假设下, 总共需要的时间很显然是 $6cn \log n$ 。那么对于稠密多项式, 我们希望 p^4 是通过两次平方的方式计算得到。

现在假设 $p(x)$ 是有 n 项的稀疏多项式。如果利用算法 8.8 按照 $(p^2)^2$ 的方式计算, 则第一次平方需要的时间是 $cn^2 \log n$, 在假设只有很少项可以合并的情况下, 第二次平方需要时间 $cn^4 \log n^2$ 。那么整个计算过程需要花费 $c(2n^4 + n^2) \log n$ 时间。利用另外一种方法, 按照 $p^4 = p \times (p \times (p \times p))$ 计算, 则总共需要 $cn^2 \log n + cn^3 \log n + cn^4 \log n = c(n^4 + n^3 + n^2) \log n$ 时间。这个时间比两次平方需要的时间更少。由此可以看出, 通过两次平方的方法并不总是计算稀疏多项式 $p^4(x)$ 的好方法。如果计算 p^{2^t} , 那么这两种方法之间的区别会更明显, 其中 t 是一个大的整数。

习题

8.1 利用算法 8.1 计算 429 的“倒数”。

8.2 利用算法 8.2 计算 429^2 。

8.3 利用算法 8.3 计算下式的“倒数”

$$x^7 - x^6 + x^5 - 3x^4 + x^3 - x^2 + 2x + 1$$

*8.4 设计一个和算法 8.2 类似的算法来计算多项式的平方。

8.5 利用习题 8.4 中编写的算法计算 $(x^3 - x^2 + x - 2)^2$ 。

8.6 利用算法 8.4 寻求 1 000 000 的表示, 其中模为 2, 3, 5, 7, 11, 13, 17, 19。

8.7 编写一个完整的算法以计算多项式模一组多项式模(a collection of polynomial moduli)的余数。

8.8 计算 $x^7 + 3x^6 + x^5 + 3x^4 + x^2 + 1$ 模 $x + 3$, $x^3 - 3x + 1$, $x^2 + x - 2$ 和 $x^2 - 1$ 的余数。

8.9 习题 8.8 中的多项式是通过仔细挑选过的, 你可以通过手工计算得到结果。随机选择次数为 1、2、3、4 的四个多项式, 并且计算 $x^9 - 4$ 除以这 4 个多项式的余数, 结果如何?

8.10 假设 5、6、7、11 是四个模, 找一个小于它们乘积的数 u , 使 $u \leftrightarrow (1, 2, 3, 4)$ 。

*8.11 推广引理 8.3, 运用到任意的多项式模上, 且这些模的根已知或者可以找到(例如次数小于 5 的多项式)。分析一下中国余数多项式在利用你编写的算法作为计算数或者模的次数的函数时的复杂度。

8.12 寻找一个多项式, 要求它在 0、1、2、3 的值分别是 1、1、2、2。

8.13 寻找下面两个多项式的最大公因子,

$$x^6 + 3x^5 + 3x^4 + x^3 - x^2 - x - 1$$

$$x^6 + 2x^5 + x^4 + 2x^3 + 2x^2 + x + 1$$

8.14 习题 8.13 中的多项式是通过挑选的, 可以通过手工计算得到答案。选择次数为 7, 和 8 的任意两个多项式, 计算它们的最大公因子。结果如何? 你怀疑得到的最大公因子吗?

*8.15 设计一个完整的算法, 用于计算整数的最大公因子, 要求该算法对 n 位整数的运行时间在 $O_b(n \log^2 n \log \log n)$ 内。

8.16 利用习题 8.15 所设计的算法计算 $\text{GCD}(377, 233)$ 。

*8.17 假设 $p(x)$ 是一个稀疏 n 次多项式。编写一个 n 的函数, 能够决定计算 $p^8(x)$ 的最好方法, 其中所用的乘法为算法 8.8。

- **8.18 下面是计算多项式 $f(x)/g(x)$ 的方法, 假设 g 能整除 f , 令 f 和 g 是次数少于或等于 $n-1$ 的多项式。
- (1) 分别计算 f 和 g 的离散傅里叶变换 F 和 G 。
 - (2) 把 F 中的项分别除以 G 中的对应项, 得到序列 H 。
 - (3) 对 H 实现逆变换, 得到结果 f/g 。这个算法是否有效。
- **8.19 假设 $M(n)$ 表示 n 位数相乘需要的时间, $Q(n)$ 表示由 n 位整数 i 计算 $\lfloor \sqrt{i} \rfloor$ 的时间。假设当 $a \geq 1$ 时, $M(an) \geq aM(n)$ 成立, 对 $Q(n)$ 也类似。证明 $M(n)$ 和 $Q(n)$ 关于常数因子是相等的。
- *8.20 把习题 8.19 推广到 (a) 多项式情况和 (b) 对固定的 r , 求 r 次方根。
- **8.21 给定计算某点 $(n-1)$ 次多项式及其该点的所有导数值的算法, 要求其复杂度为 $O_A(n \log n)$ 。
- **8.22 r 个变量[⊙]的稠密多项式可以表示为

$$\sum_{i_1=0}^{n-1} \sum_{i_2=0}^{n-1} \cdots \sum_{i_r=0}^{n-1} a_{i_1, i_2, \dots, i_r} x_1^{i_1} x_2^{i_2} \cdots x_r^{i_r}$$

证明, 在 $x_1 = \omega^{j_1}, x_2 = \omega^{j_2}, \dots, x_r = \omega^{j_r}$ 这些点上进行计值和插值操作, 可以在 $O_A(n' \log n)$ 时间内计算这类多项式的乘法, 其中 $0 \leq j_k < 2n$, ω 是主 $2n$ 次单位根。

- **8.23 证明通过对 $M(n)$ 和 $D(n)$ 平滑度上的合理假设, r 个变量的稠密多项式乘法和除法计算需要的时间 $M(n)$ 和 $D(n)$, 差别在一个常量因子之内。
- **8.24 设计一个算法能够在 $O_B(n \log^2 n \log \log n)$ 时间内, a) 把 n 位二进制数转换为十进制数; b) n 位十进制数转换为二进制数。
- **8.25 整数(多项式) x 和 y 的最小公倍数(LCM)是 z (多项式情况下可能有多个 z), 其中最小公倍数 z 能够被 x 和 y 整除, 也能够被 x 和 y 的其他整数(多项式)整除。证明计算两个 n 位数的 LCM 所需要的时间至少与两个 n 位数相乘需要的时间相同。
- 8.26 使用 2.6 节中整数乘法的方法是否能够产生 $O_A(n^{1.59})$ 时间的多项式乘法算法?
- 8.27 证明可以在 $O_A(n \log n)$ 步内对 $(n-1)$ 次多项式在点 a^0, a^1, \dots, a^{n-1} 计值。[提示: 证明存在函数 f 和 g , 可以把 $c_j = \sum_{i=0}^{n-1} b_i a^{ij}$ 表示成 $c_j = \sum_{i=0}^{n-1} f(i) g(j-i)$ 。]
- 8.28 证明可以在 $O_A(n \log n)$ 步内对 $(n-1)$ 次多项式 $ba^{2j} + ca^j + d (0 \leq j < n)$ 在 n 点上计值。
- 8.29 设计算法 8.4 和 8.5 的递归版本。

研究性问题

- 8.30 本章证明大量多项式和整数问题是

- i) 本质上具有相同的乘法复杂度, 或者
- ii) 最多差别为比乘法更复杂一个 \log 因子。

这类问题在本章的习题中有。习题 9.9 给出了另外一个组(ii)中的问题, 即“与或”乘。

一个合理的研究问题是添加问题到(i)或者(ii)中。另外一个方面就是证明一些在(ii)中的问题在本质上有相同的复杂度。例如: 可以猜想对于 GCD 和 LCM 其是正确的。

- 8.31 另外一个看似简单的问题是, 对有序稀疏多项式乘法确定算法 8.8 是否产生有序结果的最好算法。Johnson [1974] 对这个问题做了一定的研究。
- 8.32 这里描述的许多算法在实际应用中要求 n 的值非常大。然而, 对于小的 n , 可以结合 2.6 节提到的 $O(n^{1.59})$ 的方法和习题 8.26 提到的 $O(n^2)$ 的方法, 以及显而易见的 $O(n^2)$ 方法相

⊙ 随着 r 越来越大, 稠密假设就变得越来越没意义。

结合。对于小于某个上界的 n ，比本章叙述的技术有更好性能的算法。Kung [1973] 做过一些这方面的工作。

文献与注释

定理 8.2 计算倒数不会比乘法的计算更难这一事实是由 Cook [1966] 发现的。奇怪的是，类似的多项式算法在几年内没有实现。Moenck 和 Borodin [1972] 提出了一个除法算法，只需要 $O_b(n \log^2 n \log \log n)$ 时间。稍后有几个人，包括 Sieveking [1972]，独立地提出了需要 $O_b(n \log n \log \log n)$ 时间的除法算法。

模运算、插值和多项式计值算法的发展主要是 Moenck 和 Borodin [1972] 的工作。Heindel 和 Horowitz [1971] 发现了需要 $O_b(n \log^2 n \log \log n)$ 时间的基于预计算的中国余数算法。Borodin 和 Munro [1971] 提出了一个 $O(n^{1.91})$ 时间的多项式多点计值算法，Horowitz [1972] 把这个算法扩展到插值中。Kung [1973] 提出了一个 $O_b(n \log^2 n)$ 时间的多项式计值算法，不使用快速的 $(n \log n)$ 除法算法。中国余数、插值、多项式计值和整数余数的计算综合在 Lipson [1971] 中叙述。

$O_b(n \log^2 n \log \log n)$ 时间的整数最大公因子算法是由 Schönhage [1971] 提出的，其适用于多项式和一般的欧几里得域由 Moenck [1973] 提出。

经典的公因子技术的概述见 Knuth [1969]。关于稀疏多项式复杂度的材料实例由 Brown [1971]、Gentleman [1972] 和 Gentleman 和 Johnson [1973] 提出并进行研究。多项式算术的计算机实现由 Hall [1971]、Brown [1973] 和 Collins [1973] 进行。算法 8.8 的基于堆数据结构的实现由 S. Johnson 在 ALTRAN [Brown, 1973] 完成。

习题 8.19 和 8.20 由 R. Karp 和 J. Ullman 提出。习题 8.21 中的多项式计值及其导数的计算分别由 Vari [1974] 和 Aho、Steiglitz 和 Ullman [1974] 提出。习题 8.28 则稍后形成。习题 8.27 则由 Bluestein [1970] 和 Rabiner、Schafer 和 Rader [1969] 提出。多项式和整数算术处理的扩展由 Borodin 和 Munro [1975] 发现。

第9章 模式匹配算法

模式匹配是许多文本编辑、数据检索和符号操作问题的重要组成部分。在典型的字符串匹配问题中，一般会给定文本串 x 和模式串集合 $\{y_1, y_2, \dots\}$ ，要求定位模式串在 x 中的一次出现或所有的出现。模式集合通常是用正则表达式表示的正则集合。本章将给出几个解决模式匹配问题的相关技术。

本章首先回顾一下正则表达式和有穷自动机的概念，然后给出在串 x 中查询子串 y 的出现的算法，子串 y 属于给定的正则表达式所表示的集合。算法运行的时间复杂度与 x 的长度和正则表达式长度的乘积属于同一数量级。接下来给出一个线性算法用于确定字符串 y 是否为字符串 x 的子串。在这之后，给出一个功能强大的理论结果：任何能够用双向确定型下推自动机解决的模式识别问题都能够用 RAM 在线性时间内解决。因为下推自动机可能需要二次幂或甚至指数数量级的时间复杂度，所以这个结论很重要。本章最后介绍位置树(position tree)的概念，并将这个概念应用于一些模式匹配问题，如在给定串中查询最长重复子串问题。

9.1 有穷自动机和正则表达式

许多模式识别问题及其解决方法都能够用正则表达式和有穷自动机来表示，因此我们先来回顾一下相关内容。

定义 字母表是符号的有穷集合。字母表 I 上的字符串是由 I 上的符号组成的有限长度的序列。空串用 ϵ 表示，指没有任何符号的字符串。如果 x 和 y 是字符串，则 x 和 y 的连接(concatenation)是字符串 xy 。如果 xyz 是字符串，则 x 是 xyz 的前缀， y 是子串， z 是后缀。字符串 x 的长度用 $|x|$ 表示，指 x 中符号的总个数。例如，字符串 aab 的长度为 3； ϵ 的长度为 0。

字母表 I 上的一个语言是 I 上的字符串的集合。设 L_1 和 L_2 是两个语言， $L_1 L_2$ 称为语言 L_1 和 L_2 的连接，可表示为 $\{xy \mid x \in L_1, y \in L_2\}$ 。

设 L 是一个语言，定义 $L^0 = \{\epsilon\}$ ， $L^i = LL^{i-1}$ ($i \geq 1$)。 L 的克莱因闭包(Kleene closure)记作 L^* ，是语言 $L^* = \bigcup_{i=0}^{\infty} L^i$ ， L 的正闭包记作 L^+ ，是语言 $L^+ = \bigcup_{i=1}^{\infty} L^i$ 。

正则表达式和它们表示的语言(即正则集)是计算机科学的很多领域中非常有用的概念。在本章中，我们将看到它们是关于模式的有用的描述子。

定义 设 I 是一个字母表， I 上的正则表达式和其表示的语言可递归定义如下：

- 1) \emptyset 是一个正则表达式，表示空集。
- 2) ϵ 是一个正则表达式，表示集合 $\{\epsilon\}$ 。
- 3) 对 I 中的每个 a ， a 是一个正则表达式，表示集合 $\{a\}$ 。
- 4) 如果 p 和 q 是正则表达式，分别表示正则集 P 和 Q ，则 $(p+q)$ 、 (pq) 和 (P^*) 也是正则表达式，分别表示集合 $P \cup Q$ 、 PQ 或 P^* 。

在写正则表达式时，如果假设 $*$ 的运算优先级高于连接运算或 $+$ 运算，且连接运算的优先级高于 $+$ 运算，则可以省略许多括号。例如， $((0(1^*))+0)$ 可写成 01^*+0 。表达式 pp^* 可简写成 p^* 。

例 9.1

1. 01 是正则表达式, 表示集合 $\{01\}$ 。

2. $(0+1)^*$ 表示 $\{0, 1\}^*$ 。

3. $1(0+1)^*1+1$ 表示所有以 1 开头和结尾的字符串的集合。□

当且仅当一个语言能被一个正则表达式表示时, 定义该语言为正则的。如果两个正则表达式 α 和 β 表示同一集合, 则称 α 和 β 是等价的, 记作 $\alpha = \beta$ 。例如, $(0+1)^* = (0^*1^*)^*$ 。

第 4 章介绍了确定型有穷自动机的概念。这种自动机可以被看作一个设备, 包括一个“控制器”和一个输入带, 控制器总是处在有穷状态集中的一个状态下, 输入带被一个磁头从左到右扫描。确定型有穷自动机根据控制器的当前状态和输入磁头下的输入符号确定下一步的“移动操作”。每次移动后, 控制器进入一个新的状态, 且输入磁头向右移动一格。关于有穷自动机的一个重要事实是当且仅当一个语言能被有穷自动机接受时, 它才能用一个正则表达式表示。

有穷自动机的主要扩展是非确定型有穷自动机。对于每个状态和输入符号, 非确定型有穷自动机在下一步移动时有 0 个或多个选择。它也可能只改变状态, 而不移动输入磁头。

定义 一个非确定型有穷自动机 (NFA) M 是一个五元组 (S, I, δ, s_0, F) , 其中:

- 1) S 是控制器的有穷状态集。
- 2) I 是字母表, 输入符号选自该字母表。
- 3) δ 是状态转换函数, 将集合 $S \times (I \cup \{\epsilon\})$ 映射到由 S 的子集构成的集合。
- 4) S 中的 s_0 是有穷控制器的初始状态。
- 5) $F \subseteq S$ 是最终 (或接受) 状态集。

NFA M 的瞬时描述 (Instantaneous description, ID) 是一个值对 (s, w) , 其中 $s \in S$, 表示有限控制器的状态; $w \in I^*$, 表示输入串中没有使用的部分 (即输入磁头所在位置右边的字符串)。对 M 的初始 ID, 第一个元素是初始状态, 第二个元素是需要识别的字符串, 即 (s_0, w) , 其中, $w \in I^*$ 。 M 的接受 ID 形如 (s, ϵ) , 其中 $s \in F$ 。

可用 ID 的二元关系 \vdash 表示 NFA 的移动。如果 $\delta(s, a)$ 包含状态 s' , 则对于所有 I^* 中的 w , 有 $(s, aw) \vdash (s', w)$ 。注意, a 可能是 ϵ 或者是 I 中的一个符号。如果 $a = \epsilon$, 则状态转换与扫描的输入符号无关 (即不管磁头扫描的输入符号是什么, M 的状态均从 s 转到 s' 。——译者注); 如果 $a \neq \epsilon$, 则 a 必须出现在磁带的下一个输入格上, 并且输入磁头向右移动一格。

用 \vdash^* 表示 \vdash 的自反传递闭包。如果 $(s_0, w) \vdash^* (s, \epsilon)$, $s \in F$, 则称 w 可以被 M 接受。也就是说, 如果对输入串 w , 存在包含 0 个或多个移动的 ID 序列, 根据这个序列, M 能够从初始的 ID (s_0, w) 移动到一个接受 ID (s, ϵ) , 则称输入串 w 被 M 接受。 M 接受的字符串的集合称为 M 接受的语言, 记作 $L(M)$ 。

例 9.2 考虑一个非确定型有穷自动机 M , M 接受由 a, b 组成的以 aba 结尾的所有字符串, 即 $L(M) = (a+b)^*aba$ 。设 $M = (\{s_1, s_2, s_3, s_4\}, \{a, b\}, \delta, s_1, \{s_4\})$, 其中 δ 的定义如图 9-1 所示 (这里没有必要作 ϵ 转换)。

假设 M 的输入是 $ababa$ 。 M 将产生如图 9-2 所示的 ID 序列。既然 $(s_1, ababa) \vdash^* (s_4, \epsilon)$, 且 s_4 是终结状态, 则 M 接受串 $ababa$ 。□

和 NFA 关联的是一个有向图, 这个有向图能够表达了 NFA 的状态转换功能。

定义 设 $M = (S, I, \delta, s_0, F)$ 是一个 NFA, 和 M 关联的转换图 (transition diagram) 是边带标记的有向图 $G = (S, E)$ 。边及标记的集合 E 定义如下。如果对 $I \cup \{\epsilon\}$ 中某个符号 a , $\delta(s, a)$ 包含 s' , 则边 (s, s') 在 E 中; (s, s') 的标记是 $I \cup \{\epsilon\}$ 中符号 a 的集合, 使得 $\delta(s, a)$ 包含 s' 。

状态	输入		
	a	b	ϵ
s_1	$\{s_1, s_2\}$	$\{s_1\}$	\emptyset
s_2	\emptyset	$\{s_3\}$	\emptyset
s_3	$\{s_4\}$	\emptyset	$\{s_1\}$
s_4	\emptyset	\emptyset	$\{s_2\}$

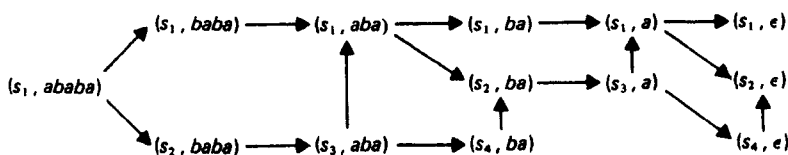
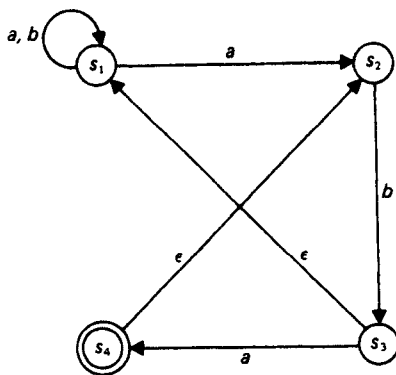
图 9-1 状态转换函数 δ 图 9-2 对应输入 $ababa$ 的 ID 序列

图 9-3 例 9.2 的转换图

例 9.3 例 9.2 的 NDFAM 的转换图如图 9-3 所示，双圈表示终结状态。 \square

用特定的闭半环将 NFA 的转换图与图中的路径问题相关联。设 I 是字母表， $S_I = (\mathcal{P}(I^*), \cup, \cdot, \emptyset, \{\epsilon\})$ 。由 5.6 节可以知道， S_I 是一个闭半环，其中 $\mathcal{P}(I^*)$ 是 I 上所有语言的集合， \emptyset 是并运算的单位， $\{\epsilon\}$ 是连接 \cdot 运算的单位。

定理 9.1 被非确定型有穷自动机接受的每个语言都是正则集。

证明： 设 $M = (S, I, \delta, s_0, F)$ 是一个 NFA， $G = (S, E)$ 是相应的转换图。对转换图中的每对顶点 s 和 s' ，应用算法 5.5 可以计算出语言 $L_{ss'}$ ，它是标记 s 到 s' 的路径的所有字符串的集合。可以看到，转换图每条边的标记都是正则表达式。此外，在算法 5.5 中，如果计算的集合 C_{ij}^{k-1} 也是正则表达式，则根据算法的第 5 行，集合 C_{ij}^k 也是正则表达式。这样，每个语言 $L_{ss'}$ 都能够用正则表达式表示，因此是一个正则集。根据定义，正则集的并也是正则的，所以 $L(M) = \bigcup_{i \in F} L_{s_0 i}$ 是一个正则集。

定理 9.1 的逆定理也成立，即给定一个正则表达式，存在一个 NFA，它可接受这个正则表达式表示的语言。从计算复杂度的角度来说，最重要的是能够找到一个 NFA，它所包含的状态数不超过正则表达式长度的两倍，且所有状态的后继状态不超过两个。

定理 9.2 如果 α 是正则表达式，则存在一个接受 α 所表示的语言的 NFA M ， $M = (S, I,$

$\delta, S_0, \{s_f\}$), 且 M 有如下属性:

1. $\|S\| \leq 2|\alpha|$, 其中 $|\alpha|$ 表示 α 的长度。[⊖]
2. 对 $I \cup \{\epsilon\}$ 中的每个字符 a , $\delta(s_f, a)$ 是空集。
3. 在每个状态 $s \in S$ 下, 对 $I \cup \{\epsilon\}$ 中所有的字符 a , $\|\delta(s, a)\|$ 的和最多为 2。

证明: 对 α 的长度进行归纳。归纳基础: 当 $|\alpha| = 1$ 时, α 必定是三种形式中的一种: (a) \emptyset , (b) ϵ 和 (c) $a (a \in I)$ 。图 9-4 所示的三个包含两个状态的自动机接受这三种形式的语言, 满足定理的条件。

归纳步骤: α 必然是以下四种形式中的一种: (a) $\beta + \gamma$, (b) $\beta\gamma$, (c) β^* 或 (d) (β) (β 和 γ 是正则表达式)。在 (d) 中, α 和 β 表示同一种语言, 所以结论显然成立。对其他情况, 设 M' 和 M'' 分别是接受 β 和 γ 所表示语言的 NDFA, 它们的状态集不相交。设它们的初始状态分别是 s_0' 和 s_0'' , 它们的终结状态分别是 s_f' 和 s_f'' 。(a)、(b)、(c) 中的转换图如图 9-5 所示。

设 α 、 β 和 γ 的长度分别为 $|\alpha|$ 、 $|\beta|$ 和 $|\gamma|$, n , n' 和 n'' 分别是 M , M' 和 M'' 中的状态数。在 (a) 中, $|\alpha| = |\beta| + |\gamma| + 1$, 且 $n = n' + n'' + 2$ 。根据归纳假设, $n' \leq 2|\beta|$, $n'' \leq 2|\gamma|$, 所以有 $n \leq 2|\alpha|$ 。并且, 在 (a) 中加入的仅有的边是从 s_0 出发的两条边, 满足从每个顶点至多引出两条边的限制, 且从 s_f' 和 s_f'' 各有一条边引出。根据假设, 之前没有从 s_f' 和 s_f'' 出发的边, 则从 s_f' 和 s_f'' 出发的边数满足边的限制。最后, 显然没有从 s_f 出发的边。同理可证明 (b)。[⊖]

对于 (c), 有 $|\alpha| = |\beta| + 1$, $n = n' + 2$, 因为 $n' \leq 2|\beta|$, 所以 $n \leq 2|\alpha|$, 很容易验证从每个顶点出发的边不超过两条的约束。

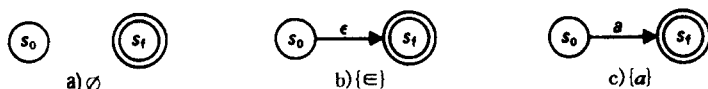


图 9-4 长为 1 的正则表达式对应的 NDFA M 的状态转换图

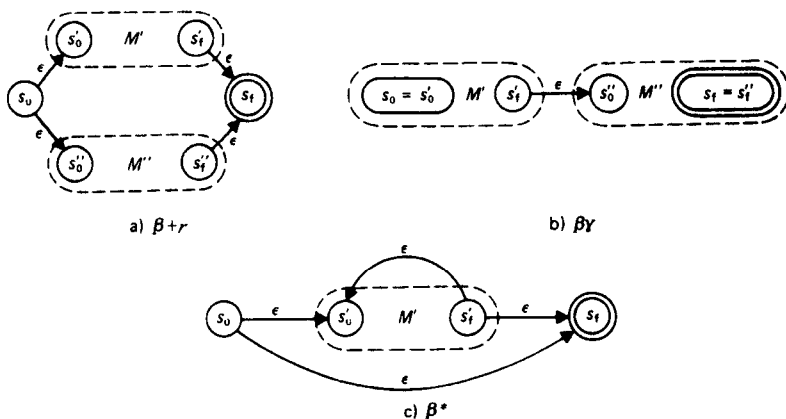


图 9-5 接受较长正则表达式表示的语言的 NDFA M 的状态转换图

- ⊖ 一个正则表达式 α 的长度是串 α 所包含的符号个数。例如, $|\epsilon^*| = 2$ 。当计算 α 的长度时, 应忽略括号 [如正则表达式 $(a^*b^*)^*$ 的“长度”是 5 而不是 7]。
- ⊖ 事实上, 从 s_f' 到 s_0'' 的边并不是必要的。我们能够识别 s_f' 和 s_0'' 。类似地, 在图 9-5a 中, 也能够识别 s_f' 和 s_f'' , 并把它们作为终结状态。

例 9.4 下面为长度为 5 的正则表达式 $ab^* + c$ 构造一个 NFA。a、b 和 c 的 NFA 如图 9-4c 所示。用图 9-5c 所示的构造方式构造 b^* 的自动机如图 9-6a 所示。用图 9-5b 所示的构造方式构造 ab^* 的自动机如图 9-6b 所示。最后，用图 9-5a 所示的构造方式构造 $ab^* + c$ 的 NFA。这个自动机有 10 个状态，如图 9-6c 所示。□

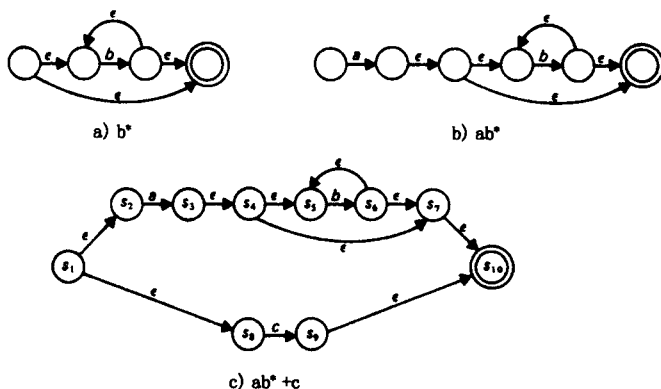


图 9-6 为 $ab^* + c$ 构造 NFA 的过程

还有另一个结果值得一提。给定任何 NFA，能够找到一个等价的“确定型”有穷自动机 DFA。然而，对于给定的 n 个状态的 NFA，等价的确定型有穷自动机的状态数有 2^n 个，因此，转换到确定型有穷自动机并不总是模拟非确定型有穷自动机的有效方法。然而，确定有穷自动机在模式识别方面是非常有用的。首先回顾一下相关的定义。

定义 确定型有穷自动机(DFA)是一个满足以下条件的非确定型有穷自动机(S, I, δ, s_0, F):

- 1) 对于所有的 $s \in S, \sigma(s, \epsilon) = \emptyset$
- 2) 对于所有的 $s \in S$ 和 $a \in I, \|\delta(s, a)\| \leq 1$ 。

定理 9.3 如果 L 是一个正则集，则存在一个 DFA 接受 L 。

证明: 由定理 9.2 知，存在一个 NFA $M = (S, I, \delta, s_0, \{s_f\})$ 接受 L ，只要将 M 转换成 DFA 即可。首先，需要找到状态对 (s, t) ，使得 $(s, \epsilon) \stackrel{*}{\Rightarrow} (t, \epsilon)$ 。为此，先构造一个有向图 $G = (S, E)$ ，图中有边 (s, t) ，(当且仅当 $\delta(s, \epsilon)$ 包含 t)，然后计算 $G' = (S, E')$ ， G' 是 G 的自反传递闭包。可以看出，当且仅当 (s, t) 在 E' 中时， $(s, \epsilon) \stackrel{*}{\Rightarrow} (t, \epsilon)$ 。

现在按照如下方式构造 NFA $M' = (S', I, \delta', s_0, F')$ ，满足 $L(M') = L(M)$ 且 M' 没有 ϵ 转换：

- 1) $S' = \{s_0\} \cup \{t \mid \delta(s, a) \text{ 包含 } t, s \in S, a \in I\}$
- 2) 对于每个 $s \in S'$ 和 $a \in I, \delta'(s, a) = \{u \mid (s, t) \in E' \text{ 且 } \delta(t, a) \text{ 包含 } u\}$
- 3) $F' = \{s \mid (s, f) \in E' \text{ 且 } f \in F\}$

证明 $L(M') = L(M)$ 的过程留给读者作为练习。这里 M' 没有 ϵ 转换。

下面，从 M' 构造一个 DFA M'' ，它的状态集是 S' 的幂集，即 $M'' = (\mathcal{P}(S'), I, \delta'', \{s_0\}, F'')$ ，其中：

- 1) 对 S' 的每个子集 S 以及 $a \in I$ ，有 $\delta''(S, a) = \{t \mid \text{对某个 } s \in S, \delta'(s, a) \text{ 包含 } t\}$
- 2) $F'' = \{S \mid S \cap F' \neq \emptyset\}$

通过对 $|w|$ 进行归纳可以证明当且仅当 $S = \{t \mid (s_0, w) \stackrel{*}{\Rightarrow} (t, \epsilon)\}, (\{s_0\}, w) \stackrel{*}{\Rightarrow} (S, \epsilon)$ ，

\in), 证明过程作为练习留给读者自行完成。因此, 可得 $L(M) = L(M') = L(M'')$ 。□

例 9.5 考虑图 9-7 的 NDFAM M , 初始状态 s_1 能够沿着标记 \in 的路径到达状态 s_3 和终结状态 s_4 , 这样, 在证明定理 9.3 时提到的计算有向图 G 的自反传递闭包 G' 时, 必须增加边 (s_1, s_4) 。完整的 G' 如图 9-8 所示。从 M 和 G' 可构造 NDFAM M' , 如图 9-9 所示。因为 G' 的每个顶点都有一条边进入 s_4 , 所以把 M' 中的所有状态都作为终结状态。既然在 M 中只有进入 s_3 的边标记为 \in , 因此 s_3 不出现在 M' 中。

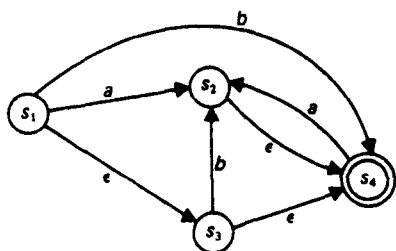


图 9-7 例 9.5 的 NDFAM

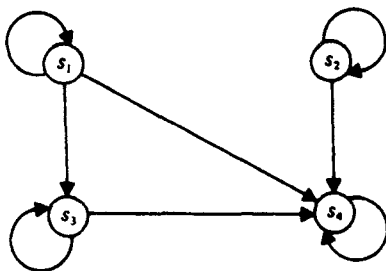


图 9-8 图 G'

在从 M' 构造 DFAM'' 的过程中, 构造了 8 个状态。然而, 只有 4 个状态是从初始状态可达的, 因此, 其他 4 个状态可以删除。最后得到的 DFAM'' 如图 9-10 所示。□

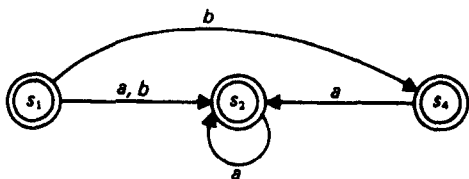


图 9-9 NDFAM M'

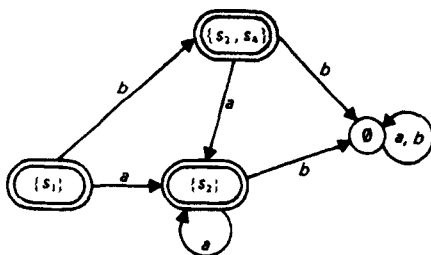


图 9-10 DFA M''

9.2 正则表达式的模式识别

考虑一个模式识别问题: 给定文本串 $x = a_1 a_2 \cdots a_n$ 和称为模式的正则表达式 α 。我们希望找到最小的 j , 对于 j , 存在某个 i , 使得 x 的子串 $a_i a_{i+1} \cdots a_j$ 属于 α 表示的语言。

这类问题在文本编辑程序中普遍存在。许多文本编辑程序允许用户在文本中进行替换操作。例如, 用户可以指出准备用某个单词代替一段文本 x 中的单词 y 。为了处理这个命令, 文本编辑程序必须能够定位作为 x 的子串的 y 在 x 中出现的位置。一些复杂的文本编辑器允许用户给出一个正则集作为要替换的可能串集合。例如, 一个用户可能说: “用空串替代 x 中的 $[I^*]$ ”, 这意味着从 x 中删除括号 $[]$ 以及位于括号中的符号。

上面问题的形式化陈述就是用表达式 $\beta = I^* \alpha$ 替换给定的正则表达式 α , 其中 I 是文本串的字母表。通过在 β 表示的语言中寻找 x 的最短前缀, 可得到 α 在 $x = a_1 a_2 \cdots a_n$ 中的第一次出现的位置。为了解决这个问题, 首先构造一个 NDFAM 来识别 β 表示的集合, 然后应用(下面的)算法 9.1, 确定状态 S_i 的集合序列。 S_i 是 NDFAM 在读入 $a_1 a_2 \cdots a_i$ ($i = 1, 2, \cdots, n$) 字符串之后所的状态集合。只

要 S_j 包含一个终结状态, 我们就知道 $a_1 a_2 \cdots a_j$ 有一个后缀 $a_i a_{i+1} \cdots a_j$, 使得 $a_i a_{i+1} \cdots a_j$ 在 α 表示的语言中(其中, $1 \leq i \leq j$)。习题 9.6 ~ 9.8 讨论了如何在字符串中找到 i (模式的最左端) 的技术。

正如定理 9.3 所述, 可以通过将 N DFA 转换成确定型有穷自动机, 从而模拟 N DFAM 对于文本串 x 的行为。然而, 这种方法的代价很大, 因为需要从正则表达式 β 得到包含 $2|\beta|$ 个状态的 N DFA, 而后再将 N DFA 转成 DFA, 而这个 DFA 的状态数将近 $2^{2|\beta|}$ 个。构造这样一个 DFA 的工作量可令人望而生畏。

另一个模拟 N DFA M 对输入串 x 的行为的方法是: 首先消除 M 中的 ϵ 转换, 再用定理 9.3 中的方法, 构造一个没有 ϵ 转换的 N DFA M' 。然后, 通过对每个 i 计算 M' 在读入 $a_1 a_2 \cdots a_i$ 之后可能的状态集 S_i ($1 \leq i \leq n$), 模拟 N DFA M' 对输入串 $x = a_1 a_2 \cdots a_n$ 的行为。事实上, 每个 S_i 是定理 9.3 中的 DFAM'' 在读入 $a_1 a_2 \cdots a_i$ 之后所处的状态。

利用这种技术, 不需要构造 M'' , 只需计算 M'' 在处理串 x 过程中出现的 M'' 的状态。为了解释怎样计算集合 S_i , 必须先说明怎样从 S_{i-1} 构造 S_i 。容易看出

$$S_i = \bigcup_{s \in S_{i-1}} \delta'(s, a_i)$$

这里 δ' 是 M' 的状态转换函数。 S_i 是最多 $2|\beta|$ 个集合的并, 每个集合至多包含 $2|\beta|$ 个成员。因为在进行集合并运算时必须消除重复成员(否则, 集合的表示可能会变得冗余), 所以对于每个输入符号, M' 显然需要做 $O(|\beta|^2)$ 步模拟操作, 或者说整个模拟需要 $O(n|\beta|^2)$ 步。

令人奇怪的是, 在许多实际情况下不用消除 ϵ 转换, 直接根据定理 9.2 从正则表达式 β 直接构造 N DFAM, 其效率反而会高。由定理 9.2 给出的关键属性是: 在转换图中, 从 M 的每个状态最多引出两条边。这个属性使我们能够证明(下面的)算法 9.1 需要 $O(n|\beta|)$ 步来模拟 β 在字符串 $x = a_1 a_2 \cdots a_n$ 上构造 N DFA 的过程。

算法 9.1 模拟一个非确定型有穷自动机。

输入: 一个 N DFA $M = (S, I, \delta, s_0, F)$ 和 I^* 中的字符串 $x = a_1 a_2 \cdots a_n$ 。

输出: 状态序列 $S_0, S_1, S_2, \dots, S_n$, 它满足

$$S_i = \{s \mid (s_0, a_1 a_2 \cdots a_i) \stackrel{*}{\vdash} (s, \epsilon)\}, 0 \leq i \leq n$$

```

1.  for  $i \leftarrow 0$  until  $n$  do
    begin
2.      if  $i = 0$  then  $S_i \leftarrow \{s_0\}$ 
3.      else  $S_i \leftarrow \bigcup_{s \in S_{i-1}} \delta(s, a_i)$ ;
        comment  $S_i$  尚未达到最终值, 现在的值相应于上面提到的集合  $T_i$ 
        的对应值;
4.      标记  $S_i$  中的每个  $t$  为 "considered";
5.      标记  $S - S_i$  中的每个  $t$  为 "unconsidered";
6.      QUEUE  $\leftarrow S_i$ ;
7.      while QUEUE 非空 do
        begin
8.            找到 QUEUE 的第一个元素  $t$  并从 QUEUE 中删除;
9.            for  $\delta(t, \epsilon)$  中的每个  $u$  do
10.               if  $u$  is "unconsidered" then
                begin
11.                    标记  $u$  "considered";
12.                    将  $u$  加到 QUEUE 和  $S_i$  上
                end
            end
        end
    end
end

```

图 9-11 模拟一个非确定型有穷自动机的算法

方法: 为了从 S_{i-1} 计算 S_i , 首先找到状态集 $T_i = \{t \mid \text{对 } s \in S_{i-1}, t \in \delta(s, a_i)\}$, 然后计算 T_i 的“闭包”, 可以通过对 T_i 中的 t , 求所有的状态 $u \in \delta(t, \epsilon)$, 并将 u 加到 T_i 中来实现闭包计算。算法中 T_i 的闭包是 S_i , 可通过将 T_i 中还没有计算过 $\delta(t, \epsilon)$ 的状态 t 组织成一个队列来计算。算法如图 9-11 所示。□

例 9.6 设 M 为图 9-6c 表示的 N DFA。假设输入是 $x = ab$, 那么在第 2 行上, $S_0 = \{s_1\}$ 。在第 9~12 行上, s_1 使 s_2 和 s_8 添加到 QUEUE 和 S_0 。考虑 s_2 和 s_8 则不会添加任何东西, 因此, $S_0 = \{s_1, s_2, s_8\}$ 。在第 3 行, $S_1 = \{s_3\}$, 考虑 s_3 导致 s_4 添加到 S_1 上, s_4 又引起 s_5 和 s_7 添加到 S_1 上, 而 s_5 什么也不添加, s_7 使 s_{10} 添加到 S_1 上, 因此, $S_1 = \{s_3, s_4, s_5, s_7, s_{10}\}$, 在第 3 行上, S_2 被设置成 $\{s_6\}$ 。 s_6 使 s_5 和 s_7 添加到 S_2 中, 而 s_7 导致 s_{10} 添加到 S_2 , 因此, $S_2 = \{s_5, s_6, s_7, s_{10}\}$ 。□

定理 9.4 算法 9.1 正确地计算状态序列 S_0, S_1, \dots, S_n , 其中 $S_i = \{s \mid (s_0, a_1 a_2 \dots a_i)^*(s, \epsilon)\}$ 。

证明: 算法 9.1 的正确性证明是一个简单的归纳证明, 留给读者作为练习。□

定理 9.5 假设自动机 M 的转换图上从每个顶点引出的边不超过 e 条, 且 M 有 m 个状态, 则算法 9.1 对于长度为 n 的输入串需要进行 $O(emn)$ 步计算。

证明: 考虑某个特定 i 值的 S_i 的计算, 图 9-11 的第 8~12 行需要执行 $O(e)$ 步。因为对于给定的 i , 没有状态会被放到 QUEUE 上两次, 所以 7~12 行的循环需要执行 $O(em)$ 时间, 容易看出, 主循环体(行 2~12)执行需要的时间为 $O(em)$ 。因此, 整个算法需要执行 $O(emn)$ 步。

下面是一个重要的推论, 该推论将正则集的识别与算法 9.1 相关联。

推论 如果 β 是正则表达式, $x = a_1 a_2 \dots a_n$ 是长度为 n 的字符串, 那么有一个接受 β 表示的语言 N DFA M , 使得算法 9.1 需要 $O(n|\beta|)$ 步确定状态序列 S_0, S_1, \dots, S_n , 这里 $S_i = \{s \mid (s_0, a_1 a_2 \dots a_i)^*(s, \epsilon)\}, 0 \leq i \leq n$ 。

证明: 通过定理 9.2, 可以构造 N DFAM, M 最多有 $2|\beta|$ 个状态, 每个状态最多引出 2 条边, 这样定理 9.5 的 e 最多为 2, 根据定理 9.5, 结论显然成立。□

各种模式识别算法可以从算法 9.1 构造。例如, 假设给定正则表达式 α 和文本串 $x = a_1 a_2 \dots a_n$, 我们希望找到最小的 k , 有 $j < k$, 使 $a_j a_{j+1} \dots a_k$ 在 α 表示的集合中。应用定理 9.2, 我们能够从 α 构造一个 N DFA M 来接受语言 $I^* \alpha$ 。为了找到最小的 k , 使得 $a_1 a_2 \dots a_k$ 在 $L(M)$ 中, 我们可以在图 9-11 中的算法的第 2~12 行的语句块结束处插入一个测试语句, 用于检测 S_i 是否包含 F 中的一个状态。根据定理 9.2, 我们可以将 F 取为一个单元集合, 所以这个测试不需要花很多时间, 时间复杂度为 $O(m)$, 其中 m 是 M 的状态数。如果 S_i 包含 F 中的一个状态, 找到 $L(M)$ 中 x 的最短前缀 $a_1 a_2 \dots a_i$, 算法就可以立即跳出主循环。

可以进一步修改算法 9.1 从而为每个 k 产生最大的(或者最小的 j) $j < k$, 使 $a_j a_{j+1} \dots a_k$ 在 α 表示的集合中, 这可以通过将集合 S_i 中的每个状态与一个整数相关联予以实现。和 S_k 中的状态 s 关联的整数表示最大(或最小)的 j , 使得 $(s_0, a_j a_{j+1} \dots a_k)^*(s, \epsilon)$ 。将算法 9.1 中这些整数更新的具体细节留给读者自行完成。

9.3 子串识别

上节所述一般问题中的一个重要特例是, 当正则表达式 α 可以表示成形式 $y_1 + y_2 + \dots + y_k$ 时, 其中每个 y_i 是某个字母表 I 上的字符串。定理 9.5 的推论意味着能在文本串 $x = a_1 a_2 \dots a_n$ 中用 $O(\ln)$ 步找到模式 y_i 的第一次出现, 这里 l 是所有 y_i 的长度总和。然而, 还可以找到一个需要 $O(l+n)$ 步的解决方案。先考虑只有一个模式串 $y = b_1 b_2 \dots b_l$ 的情况, 这里的每个 b_i 均是 I 中的符号。

从模式 y 构造一个确定型模式匹配机器 M_y , M_y 用来识别 I^*y 中的最短实例串。为了构造 M_y , 首先构造一个框架 (skeletal) DFA, 这个 DFA 包含 $l+1$ 个状态, 分别标记为 $0, 1, \dots, l$, 且在输入字符 b_i 作用下将从状态 $i-1$ 转到状态 i , 如图 9-12 所示。对所有的输入 $b \neq b_1$, 状态 0 均转到它自身。可以把状态 i 看成一个指针, 指向模式串 y 中的第 i 个位置。

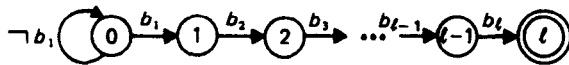


图 9-12 一个框架机器

模式匹配机器 M_y 的操作类似确定型有穷自动机, 不同之处仅在于它在扫描同样的输入字符时能够进行多种状态转换。 M_y 的状态集和框架 DFA 一样, 这样 M_y 的状态 j 对应于模式串 y 的前缀 $b_1b_2 \dots b_j$ 。

M_y 从状态 0 开始, 且输入指针指向 a_1 , a_1 是文本串 $x = a_1a_2 \dots a_n$ 的第一个字符。如果 $a_1 = b_1$, 则 M_y 进入状态 1, 并且将它的输入指针前移到文本串的位置 2。如果 $a_1 \neq b_1$, 则 M_y 仍旧在状态 0, 并且将它的输入指针前移到位置 2。

假设读入字符 $a_1a_2 \dots a_k$ 后, M_y 在状态 j , 这意味着 $a_1a_2 \dots a_k$ 的最后 j 个字符是 $b_1b_2 \dots b_j$, 且对于 $m > j$, $a_1a_2 \dots a_k$ 的最后 m 个字符不是 $b_1b_2 \dots b_l$ 的前缀。如果下一个输入字符 a_{k+1} 和 b_{j+1} 一致, M_y 进入状态 $j+1$, 且输入指针前移到 a_{k+2} ; 如果 $a_{k+1} \neq b_{j+1}$, M_y 进入编号最高的状态 i , $b_1b_2 \dots b_i$ 是 $a_1a_2 \dots a_k a_{k+1}$ 的一个后缀。

为了帮助确定状态 i , 机器 M_y 用一个整数值函数 f 和状态关联, 称为失败函数 (failure function), 使得 $f(j)$ 是小于 j 的最大的 s , ($b_1b_2 \dots b_s$ 是 $b_1b_2 \dots b_j$ 的一个后缀), 即 $f(j)$ 是最大的 $s < j$, 使得 $b_1b_2 \dots b_s = b_{j-s+1}b_{j-s+2} \dots b_j$ 。如果没有这样的 $s \geq 1$, 则 $f(j) = 0$ 。

例 9.7 假设 $y = aabbaab$, f 的值如下:

i	1	2	3	4	5	6	7
$f(i)$	0	1	0	0	1	2	3

例如, 因为在 $aabbaa$ 的所有前缀中, 能够作为 $aabbaa$ 后缀的最长的适当前缀是 aa , 所以 $f(6) = 2$ 。□

接下来给出一个算法来计算失败函数。首先, 为了说明 M_y 怎样使用失败函数, 定义函数 $f^{(m)}(j)$ 如下:

$$1) f^{(1)}(j) = f(j)$$

$$2) f^{(m)}(j) = f(f^{(m-1)}(j)), \text{ 对于 } m > 1$$

即 $f^{(m)}(j)$ 只是将 f 运用到 j 上 m 次。[在例 9.7 中, $f^{(2)}(6) = 1$ 。]

假设 M_y 读入字符串 $a_1a_2 \dots a_k$ 后进入状态 j , 且 $a_{k+1} \neq b_{j+1}$ 。此时, M_y 重复地应用失败函数到 j 上, 直到找到 m 的最小值, 使得

$$1) f^{(m)}(j) = u \text{ 且 } a_{k+1} = b_{u+1}$$

$$2) f^{(m)}(j) = 0 \text{ 且 } a_{k+1} \neq b_1$$

即 M_y 通过状态 $f^{(1)}(j)$, $f^{(2)}(j)$, \dots 回退直到情况 1 或情况 2 对 $f^{(m)}(j)$ 成立, 但是对 $f^{(m-1)}(j)$ 不成立为止。如果情况 1 成立, 则 M_y 进入状态 $u+1$ 。如果情况 2 成立, 则 M_y 进入状态 0。在任一种情况下, 输入指针都移动到位置 a_{k+2} 。

在情况 1 中, 很容易验证, 如果 $b_1b_2 \dots b_j$ 是 y 中能作 $a_1a_2 \dots a_k$ 后缀的最长前缀, 则 $b_1b_2 \dots b_{f^{(m)}(j)+1}$ 是 y 中能作为 $a_1a_2 \dots a_{k+1}$ 后缀的最长前缀。在第 2 种情况下, y 没有前缀是 $a_1a_2 \dots a_{k+1}$ 的后缀。

接下来, M_y 处理输入符号 a_{k+2} 。 M_y 用这种方式继续操作, 直到获得以下两种结果之一为止。一种是 M_y 进入终结状态 l , 此时被扫描到的最后 l 个输入符号构成模式 $y = b_1 b_2 \cdots b_l$ 的一个实例; 另一种结果是直到 M_y 处理了 x 的最后一个输入符号, 也没能进入状态 l , 此时我们知道 y 不是 x 的子串。

例 9.8 设 $y = aabbaab$ 。一个模式匹配机 M_y 如图 9-13 所示, 虚线箭头指向失败函数在每个状态下的值。在输入 $x = abaabaabbaab$ 时, M_y 将执行如下的状态转换序列:

输入: $a \ b \ a \ a \ b \ a \ a \ b \ b \ a \ a \ b$
 状态: $0 \ 1 \ 0 \ 1 \ 2 \ 3 \ 1 \ 2 \ 3 \ 4 \ 5 \ 6 \ 7$
 0 0

例如, 初始的 M_y 处在状态 0, 在读入 x 的第 1 个符号后, M_y 进入状态 1。因为在状态 1 下输入第 2 个符号 b 没有相应的状态转换, 所以 M_y 进入状态 0, 即是状态 1 下的失败函数的值, 且不动输入指针。因为 y 的第一个符号不是 b , 属于情况 2, 则 M_y 仍旧在状态 0, 继续移动它的输入指针到位置 3。

在读入第 12 个符号时, M_y 进入终结状态 7。这样, 在 x 的第 12 个位置上, M_y 找到了模式 y 的一次出现。 □

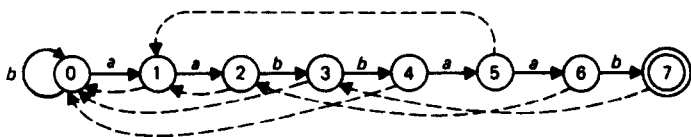


图 9-13 一个模式匹配机器

函数 f 能够以和 M_y 操作非常类似的方式进行迭代计算。根据定义, $f(1) = 0$ 。假设已经计算了 $f(1), f(2), \dots, f(j)$, 且设 $f(j) = i$ 。为了计算 $f(j+1)$, 可检查 b_{j+1} 和 b_{i+1} 。如果 $b_{j+1} = b_{i+1}$, 则有 $b_1 b_2 \cdots b_i b_{i+1} = b_{j+1} b_{j+2} \cdots b_j b_{j+1}$, 因此 $f(j+1) = f(j) + 1$; 如果 $b_{j+1} \neq b_{i+1}$, 我们找到最小的 m , 使得

$$1. f^{(m)}(j) = u, \text{ 且 } b_{j+1} = b_{u+1}$$

$$2. f^{(m)}(j) = 0 \text{ 且 } b_{j+1} \neq b_1$$

在第 1 种情况下, 令 $f(j+1) = u + 1$; 对于第二种情况, 令 $f(j+1) = 0$ 。下面算法给出相关的细节。

算法 9.2 计算失败函数。

输入: 模式 $y = b_1 b_2 \cdots b_l, l \geq 1$ 。

输出: y 的失败函数 f 。

方法: 执行如图 9-14 所示的程序。 □

例 9.9 考虑算法 9.2 在输入 $y = aabbaab$ 时的行为。开始时有 $f(1) = 0$, 因为 $b_2 = b_1$, 所以, $f(2) = 1$ 。然而, $b_3 \neq b_2$ 且 $b_3 \neq b_1$, 所以 $f(3) = 0$ 。继续这种方式, 可以得到例 9.7 中给出的 f 的函数值。 □

现在证明算法 9.2 在 $O(|y|)$ 时间内正确计算函数 f 。首先证明算法 9.2 的正确性。

定理 9.6 算法 9.2 计算失败函数 f 。

证明: 通过在 j 上应用数学归纳法来证明定理, 若 $f(j)$ 是小于 j 的最大整数 i , 使得 $b_1 b_2 \cdots b_i = b_{j-i+1} b_{j-i+2} \cdots b_j$, 如果没有这样的 i 存在, 则 $f(j) = 0$ 。

根据定义, $f(1) = 0$ 。假设对所有的 $f(k) (k < j)$, 归纳假设为真。在计算 $f(j)$ 时, 算法 9.2

在第4行比较 $b_{f(j-1)+1}$ 和 b_j 。

情况1 假设 $b_j = b_{f(j-1)+1}$ 。由于 $f(j-1)$ 是最大的 i ，使得 $b_1 \cdots b_i = b_{j-i} \cdots b_{j-1}$ ，所以 $f(j) = i + 1$ 。这样，根据算法第5行和第6行，算法正确地计算了 $f(j)$ 。

情况2 假设 $b_j \neq b_{f(j-1)+1}$ ，那么必须找到满足 $b_1 \cdots b_i = b_{j-i} \cdots b_{j-1}$ 和 $b_{i+1} = b_j$ 的 i 最大的值（如果 i 存在的话）；如果不存在这样的 i ，则很明显 $f(j) = 0$ ，而算法在第5行中已正确计算了 $f(j)$ 。假设 i_1, i_2, \dots 分别是使下式成立的 i 值的最大者，次大者，……

```

begin
1.   f(1) ← 0;
2.   for j ← 2 until l do
      begin
3.     i ← f(j-1);
4.     while b_j ≠ b_{i+1} and i > 0 do i ← f(i);
5.     if b_j ≠ b_{i+1} and i = 0 then f(j) ← 0
6.     else f(j) ← i + 1
      end
end

```

图 9-14 失败函数的计算

$$b_1 b_2 \cdots b_i = b_{j-i} \cdots b_{j-1}$$

根据简单归纳假设，有 $i_1 = f(j-1)$ ， $i_2 = f(i_1) = f^{(2)}(j-1)$ ， \dots ， $i_k = f(i_{k-1}) = f^{(k)}(j-1)$ ，因为 i_{k-1} 是第 $k-1$ 大的 i 的值，使得 $b_1 \cdots b_i = b_{j-i} \cdots b_{j-1}$ ，且 i_k 是小于 i_{k-1} 的最大的 i 值，使得 $b_1 \cdots b_i = b_{i_{k-1}-i+1} \cdots b_{i_{k-1}} = b_{j-i} \cdots b_{j-1}$ 。算法第4行依次计算比较 i_1, i_2, \dots ，直到找到 i 值（如果存在这样的 i ），满足 $b_1 \cdots b_i = b_{j-i} \cdots b_{j-1}$ 且 $b_{i+1} = b_j$ 为止。在 **while** 语句执行终止的时候，如果有存在这样一个 i_m ，则 $i = i_m$ 。由此，算法在第5行语句正确计算了 $f(j)$ 。

这样，对所有的 j ，算法都正确地计算 $f(j)$ ，定理得以证明。 \square

定理 9.7 算法 9.2 可以在 $O(l)$ 步内正确计算失败函数 f 。

证明：算法的第3行和第5行的代价是固定的，**while** 语句的代价是和第4行上 **do** 语句后的语句 $i \leftarrow f(i)$ 产生的递减次数 i 成比例的。增加 i 值的唯一方法是在第6行赋值 $f(j) = i + 1$ ，然后在第2行给 j 加1，在第3行设置 i 值为 $f(j-1)$ 。既然 i 的初始值为0，那么第6行至多执行 $l-1$ 次。于是得出结论：第4行的 **while** 语句的执行次数不可能多于 l 次。这样，执行第4行的总代价是 $O(l)$ 。算法的余下部分的代价显然是 $O(l)$ ，这样，算法总的执行时间代价为 $O(l)$ 。 \square

与定理 9.6 一样，我们能够证明当且仅当 $b_1 b_2 \cdots b_i$ 是作为 $a_1 a_2 \cdots a_k$ 的后缀的 y 的最长前缀时，模式匹配机 M_y 在读入 $a_1 a_2 \cdots a_k$ 后进入状态 i 。这样， M_y 就能够正确地在文本串 $x = a_1 a_2 \cdots a_n$ 的最左边找到 y 的出现。

根据定理 9.7 可以证明， M_y 在处理输入串 x 时，将经过至多 $2|x|$ 次状态转换。这样，我们能够通过跟踪 M_y 在输入 x 上的状态转换确定 y 是否是 x 的一个子串^①。为此，只需要 y 的失败函数 f 就可以进行判断。由定理 9.7 可知，能够在 $O(|y|)$ 时间内构造函数 f ，这样能够在 $O(|x| + |y|)$ 时间内确定 y 是否 x 的子串，这个时间与字母表的规模无关。如果模式串的字母表比较小，且文本串又比模式串长得多，那么可以考虑模拟 DFA 接受语言 I^*y 的过程来判断 y 是否是 x 的子串。该 DFA 对于输入的每个字符恰好做一次状态转换。

算法 9.3 为 I^*y 构造一个 DFA。

① 回想一下，事实上， M_y 的状态由指向模式 y 中的位置指针反映出来。因此， M_y 的状态转换能够通过 y 内移动指针直接实现。

输入：字母表 I 上的模式串 $y = b_1 b_2 \cdots b_l$ 。为了方便起见，取 b_{l+1} 为一个新的、和 I 上的任何字符都不相同的字符。

输出：一个 DFA M ，满足 $L(M) = I^* y$ 。

方法：

1. 用算法 9.2 为 y 构造失败函数 f 。
2. 设 $M = (S, I, \delta, 0, \{l\})$ ， $S = \{0, 1, \dots, l\}$ ，构造 δ 的过程如下：

```

begin
  for  $j = 1$  until  $l$  do  $\delta(j-1, b_j) \leftarrow j$ ;
  for each  $b$  in  $I$ ,  $b \neq b_1$  do  $\delta(0, b) \leftarrow 0$ ;
  for  $j = 1$  until  $l$  do
    for each  $b$  in  $I$ ,  $b \neq b_{j+1}$  do  $\delta(j, b) \leftarrow \delta(f(j), b)$ 
end

```

□

定理 9.8 算法 9.3 构造一个 DFA M ，满足

$$(0, a_1 a_2 \cdots a_k) \models^* (j, \epsilon)$$

当且仅当 $b_1 b_2 \cdots b_j$ 是 $a_1 a_2 \cdots a_k$ 的一个后缀，且不存在 $i > j$ 时， $b_1 b_2 \cdots b_i$ 是 $a_1 a_2 \cdots a_k$ 的一个后缀。

证明：利用定理 9.6 中的论证方法，对 k 进行归纳。具体证明过程留给读者自行完成。 □

例 9-10 使用算法 9.3 得出的识别模式 $y = aabbaab$ 的 DFA M 如图 9-15 所示。

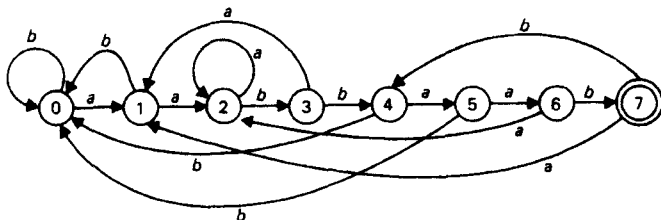


图 9-15 一个接受 $(a+b)^* aabbaab$ 的确定型有穷自动机

在输入 $x = abaabaabbaab$ 时， M 将经历以下的状态转换：

输入： a b a a b a a b b a a b

状态： 0 1 0 1 2 3 1 2 3 4 5 6 7

M 和 M_j 之间唯一的区别是 M 在字符不匹配模式串的情况下提前计算了 M 的下一个状态。这样， M 在每个输入符号上恰好做一个状态转换。

总结本节的主要结论可得出下面的定理。

定理 9.9 在 $O(|x| + |y|)$ 时间内能够确定 y 是否是 x 的子串。

现在考虑给定若干个模式串 y_1, y_2, \dots, y_k 时，确定某个 y_i 是否给定字符串 $x = a_1 a_2 \cdots a_n$ 的子串。本节介绍的方法也适用于解决这个问题。首先为 y_1, y_2, \dots, y_k 构造一个框架机器，这个框架机器将是一棵树。接着在与 $l = |y_1| + |y_2| + \cdots + |y_k|$ 值成比例的时间内计算关于这棵树的失败函数。然后可以用前面所介绍的方式构造模式匹配机。这样，在 $O(l+n)$ 步内，就能够确定某个 y_i 是否为 x 的子串，具体细节留作练习请读者自行完成。

9.4 双向确定型下推自动机

一旦有人猜测有 $O(|x| + |y|)$ 的算法可用于确定 y 是否为 x 的子串，那么这样的算法就

不难找到了。然而是什么原因会让人们猜想有这样的算法存在呢？这种猜想可能缘自对双向确定型下推自动机(two-way deterministic pushdown automata, 简称 2DPDA)的研究。

2DPDA 是一种用于接受语言的特殊的图灵机。我们可以根据语言识别问题来重新形式化许多模式识别问题。例如, 如果用 L 表示语言 $\{xcy \mid x \text{ 和 } y \text{ 均在 } I^* \text{ 中}, c \notin I, \text{ 且 } y \text{ 是 } x \text{ 的子串}\}$, 那么确定 y 是否是 x 的子串等价于确定字符串 xcy 是否语言 L 中的语句。本节将说明存在一个能够识别 L 的 2DPDA。尽管 2DPDA 识别这个语言可能需要 $O(n^2)$ 时间, 但是有一种功能强大的模拟技术, 采用这种技术, 一个 RAM 能够在 $O(n)$ 时间内模拟给定的 2DPDA 在长度为 n 的输入串上的行为。本节将详细研究这种模拟技术。

可以将一个 2DPDA 看作一个双带图灵机, 在这个图灵机上, 一个磁带用作下推存储器, 如图 9-16 所示。2DPDA 有一个只读输入带, 其最左边格子有一个左终端标记符 ϵ , 最右端格子有一个右终端标记符 $\$$, 两个终端标记符之间的是要识别的输入串, 每个符号占一格。输入符号来自输入字母表 I , 假设 I 不包含 ϵ 和 $\$$ 。输入磁头一次读一个符号, 在一次移动中输入磁头可以向左移动一格、保持静止或者向右移动一格。假设磁头不可能移到输入磁带两端的终端标记符之外; 由于存在 ϵ 和 $\$$ 标记符, 我们可以写一个移动函数(next-move), 这个函数从来不会从 ϵ 向左, 或从 $\$$ 向右移动输入磁头。

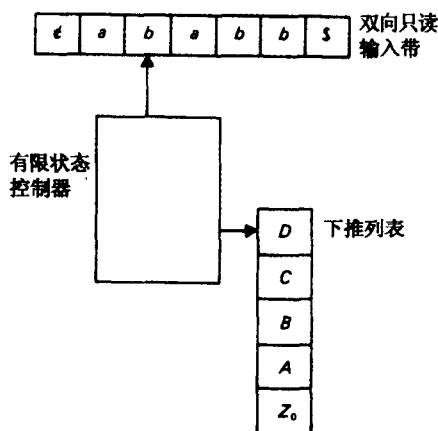


图 9-16 双向确定型下推自动机

下推列表(pushdown list)是一个栈, 用于存放来自下推列表字母表 T 中的字符。下推列表的底部单元存放特殊字符 Z_0 , 该字符标记栈底, 假设 Z_0 不在 T 中。

有限控制器总是处于有穷状态集 S 的某个状态下。机器的操作由移动函数 δ 确定, 对于 S 中的 s 、集合 $I \cup \{\epsilon, \$\}$ 中的 a 、 $T \cup \{Z_0\}$ 中的 A , $\delta(s, a, A)$ 表明当有限控制器处在状态 s 时, 输入磁头正在扫描符号 a , 下推列表中的栈顶符号为 A 。机器要采取的动作, 有三种:

$$\delta(s, a, A) = \begin{cases} (s', d, \text{push } B) & \text{已知 } B \neq Z_0 \\ (s', d) & \\ (s', d, \text{pop}) & \text{如果 } A \neq Z_0 \end{cases}$$

在这三个动作中, 机器进入状态 s' , 并沿 d 方向移动输入磁头($d = -1, +1$ 或 0 分别表示向左移动一格, 向右移动一格或保持静止)。 $\text{push } B$ 表示将符号 B 压入下推列表的顶端, pop 表示从下推列表中移出处于顶端的符号。

定义 可以形式化地将 2DPDA 定义为 7 元组:

$$P = (S, I, T, \delta, s_0, Z_0, s_f)$$

其中:

- 1) S 是有穷控制状态的集合。
- 2) I 是输入字母表(不包括 ϵ 和 $\$$)。
- 3) T 是下推列表字母表(不包括 Z_0)。

4) δ 是集合 $(S - \{s_f\}) \times (I \cup \{\epsilon, \$\}) \times (T \cup \{Z_0\})$ 上的一个映射, $\delta(s, a, A)$ 的值如果有定义, 是以下三种形式中的一种: $(s', d, \text{push } B)$, (s', d) 或 (s', d, pop) , 其中 $s' \in S$, $B \in T$ 且 $d \in \{-1, 0, +1\}$ 。假设 2DPDA 到达最终状态 s_f 后不再移动, 某些状态也不定义移动, 且对于任意的 s 和 A , $\delta(s, \epsilon, A)$ 的第二个元素 d 不是 -1 , $\delta(s, \$, A)$ 的第二个元素 d 不是 $+$ 1。最后, $\delta(s, a, Z_0)$ 没有第三个元素 pop 。

- 5) 在 S 中的 s_0 是有限控制器的初始状态。
- 6) Z_0 是下推列表的底部标记。
- 7) s_f 是指定的终结状态。

一个 2DPDA P 在输入 $w = a_1 a_2 \cdots a_n$ 上的一个瞬时描述是一个三元组 (s, i, α) , 其中:

- 1) s 是 S 中的一个状态。
- 2) i 是整数, 表示输入磁头的位置 ($0 \leq i \leq n+1$)。通常假设 $a_0 = \epsilon$, $a_{n+1} = \$$ 。
- 3) α 表示下推列表中的字符串内容, α 的最左端符号在栈顶。

将 2DPDA 在输入 $a_1 a_2 \cdots a_n$ 下的一次移动可定义成 ID 的二元关系 \vdash , 即

- 1) 如果 $\delta(s, a_i, A) = (s', d, \text{push } B)$, 则 $(s, i, A\alpha) \vdash (s', i+d, BA\alpha)$ 。
- 2) 如果 $\delta(s, a_i, A) = (s', d)$, 则 $(s, i, A\alpha) \vdash (s', i+d, A\alpha)$ 。
- 3) 如果 $\delta(s, a_i, A) = (s', d, \text{pop})$, 则 $(s, i, A\alpha) \vdash (s', i+d, \alpha)$ 。

注意, 符号只能加到下推列表的顶端, 或从顶端读取, 或从顶端移走。我们还要限制在任何时候 2DPDA 的下推列表只有一个底部标记。可以使用 \vdash^* 表示 0 个或多个移动的序列, 其中 \vdash^* 是 \vdash 的自反传递闭包。

输入 $w = a_1 a_2 \cdots a_n$ (不包括底部标记), 2DPDA P 的初始 ID 是 $(s_0, 1, Z_0)$, 表示 P 处在初始状态 s_0 , 输入磁头正在扫描 $a_1 a_2 \cdots a_n a_{n+1}$ 最左端的符号^①, 且下推列表只包含底部标记 Z_0 。

输入 $w = a_1 a_2 \cdots a_n$ 的接受 ID 的一种形式为 (s_f, i, Z_0) ($0 \leq i \leq n+1$), 表示 P 已经进入最终状态 s_f , 下推列表只包含底部标记。

如果 2DPDA P 在输入 w 时, 对于 $0 \leq i \leq |w| + 1$, 满足 $(s_0, 1, Z_0) \vdash^* (s_f, i, Z_0)$, 则称 2DPDA P 接受输入串 w 。 P 接受的所有字符串的集合记作 $L(P)$, 称为 P 接受的语言。

下面考虑一些 2DPDA 和它们所接受的语言的例子。这里不采用形式化的 7 元组, 而是通过刻画它们的行为来非形式化地描述 2DPDA。

例 9.11 考虑语言 $L = \{xyc \mid x \text{ 和 } y \text{ 在 } \{a, b\}^* \text{ 中, 且 } y \text{ 是 } x \text{ 的子串}\}$ 。2DPDA P 能够按如下步骤识别 L 。假设 P 的输入串 w 的形式为 xyc , 其中 $x = a_1 a_2 \cdots a_n$ 且 $a_i \in \{a, b\}$, $1 \leq i \leq n$ 。

1) P 向右移动输入磁头直到遇到 c 为止。

2) P 然后向左移动输入磁头, 复制 $a_n a_{n-1} \cdots a_1$ 到下推列表, 直到输入磁头到达左侧的终端标记为止。此时, P 的下推列表中有 $xZ_0 = a_1 a_2 \cdots a_n Z_0$ (xZ_0 的最左边的符号 a_1 在栈顶)。

3) P 移动输入磁头到 c 右边的第一个符号 (即 y 的第一个字符), 不改变下推列表, 并准备用 y 来匹配下推列表。

4) **while** 下推列表的顶端符号匹配输入磁头扫描到的符号 **do**
begin

① 如果 $n=0$, 即 $w = \epsilon$, 那么 P 正在扫描 $a_{n+1} = \$$ (最右端的标记符)。

从下推列表将顶端符号弹出栈;

输入磁头向右移动一格

end

5) 此时有两种可能:

a) 输入磁头已经到达 \$, 即右端的终端标记符, 此时 P 接受输入。

b) 下推列表的顶端的符号和当前的输入符号不匹配。既然之前串 y 和下推列表一直是匹配的, 因此 P 可以恢复其下推列表, 方式是通过向左移动输入磁头, 复制扫描到的输入到下推列表直到遇到输入符号 c 。这时, P 的下推列表的内容为 $a_i a_{i+1} \cdots a_n Z_0$, $1 \leq i \leq n^\ominus$ 。如果 $i = n$, 则 P 在非最终状态下停机。否则, P 从下推列表中 a_i 出栈, 并向右移动输入磁头到 y 的 \$ 的第一个符号。然后 P 回到步骤 4, 尝试找到字符串 $a_{i+1} \cdots a_n$ 的前缀 y 。

可以看到, P 用一种自然的方式来确定 y 是否是 $x = a_1 a_2 \cdots a_n$ 的子串, P 轮流对 $i = 1, 2, \dots, n$, 尝试用 $a_i a_{i+1} \cdots a_n$ 的前缀来匹配 y 。如果找到匹配, P 接受 y 。否则, P 拒绝 y 。这个过程可能要求 P 作 $O(|x| \cdot |y|)$ 次移动。□

例 9.12 考虑语言 $L = \{w \mid w \in \{a, b, c\}^*, w = xy \text{ 满足 } |x| \geq 2 \text{ 且 } x^R = x\}$ (即 w 的某个前缀是长度为 2 或更长的回文), 上标 R 表示字符串反转, 如 $(abb)^R = bba$ 。因为 $\{a, b, c\}^*$ 中的每个字符串均以单个字符构成的回文 a, b 或 c 开头, 所以需要施加 $|x| \geq 2$ 的要求。下面考虑给定输入串 $a_1 a_2 \cdots a_n$ 时, 2DPDA P 的行为:

1) P 移动输入磁头到右侧终端标记符, 同时复制输入到下推列表。这样, 下推列表包含 $a_n a_{n-1} \cdots a_2 a_1 Z_0$ 。如果 $n < 2$, 则 P 立即拒绝输入。

2) P 移动输入磁头到左侧终端标记符, 但不改变下推列表。 P 将输入磁头定位在左侧终端标记符右边的第一个符号上, 即指向 a_1 。

3) while 下推列表的顶端符号与输入磁头扫描的符号匹配 do

begin

将下推列表的顶端符号弹出栈;

输入磁头向右移动一格

end

4) 此时有两种可能:

a) 下推列表只包含 Z_0 , 此时 P 停机, 接受输入。

b) 下推列表顶端符号和当前输入符号不一致。此时 P 向左移动输入磁头, 将所读到的字符复制回下推列表的相应位置, 直到遇到左侧终端标记符为止。如果初始的输入串是 $a_1 a_2 \cdots a_n$, 则此时对某个 i , P 的下推列表中有 $a_i \cdots a_2 a_1 Z_0$ 。如果 $i = 2$, P 停机且拒绝。否则, P 从下推列表中弹出 a_i , 且输入磁头向右移动一格, 定位在左侧终端标记符右边的第一个符号处, 然后 P 再回到步骤 3。

这样, P 通过对每个 i 尝试匹配 $a_i \cdots a_2 a_1$ 和 $a_1 a_2 \cdots a_n$ 的前缀 ($2 \leq i \leq n$), 确定输入串 $a_1 a_2 \cdots a_n$ 是否以回文开头。这个过程可能需要 P 做 $O(n^2)$ 次移动。□

现在证明算法能够在 $O(|w|)$ 时间内确定 2DPDA P 是否接受输入串 w , 而无需考虑 P 事实上做多少次移动操作。在下文中, 假设围绕固定的 2DPDA $P = (S, I, T, \delta, s_0, Z_0, s_f)$ 和固定长度 n 的输入字符串 w 进行讨论。

定义 P 的一个表面格局 (surface configuration, 简称格局) 是一个三元组 (s, i, A) , 其中, $s \in S$, i 是输入磁头的位置, $0 \leq i \leq n+1$, A 是 $T \cup \{Z_0\}$ 中的一个下推列表符号。

⊖ 算法第一次运行到这里时, $i = 1$, 但后面每次运行到此时, i 将递增 1。

注意, 尽管并非每个 ID 都是一个格局, 但一个格局是一个 ID。从这个意义上说, 每个表面格局可表示多个 ID, 即表示那些下推列表“表面”(即栈顶)和格局中的符号 A 相匹配的 ID。如果存在 P 的移动的序列, 使得对于 $m \geq 0$, 有

$$\begin{aligned} C &\vdash (s_1, i_1, \alpha_1) \\ &\vdash (s_2, i_2, \alpha_2) \\ &\vdots \\ &\vdots \\ &\vdots \\ &\vdash (s_m, i_m, \alpha_m) \\ &\vdash C' \end{aligned}$$

我们说格局 $C = (s, i, A)$ 导出格局 $C' = (s', i', A)$, 记作 $C \Rightarrow C'$ 。其中, 对每个 i , 有 $|\alpha_i| \geq 2^\circ$ 。在这个移动序列里, 中间 ID 的 (s_j, i_j, α_j) 至少有两个符号在下推列表中。因为不同的格局只有 $O(n)$ 个, 而用到的不同的 ID 数量可能是指数级的, 所以我们处理的是格局而不是 ID。

定义 如果对于格局 (s, i, A) , 其 $\delta(s, a_i, A)$ 没有定义或者形式为 (s', d, pop) , 则称格局 (s, i, A) 是最终的 (terminal)。也就是说, 一个最终的格局将引起 P 停机或者从下推列表中弹出一个符号。格局 C 唯一的最终格局 C' (如果存在) 满足 $C \Rightarrow C'$, 这里 \Rightarrow 表示 \Rightarrow 的自反传递闭包, 则称 C' 是格局 C 的终结者 (terminator)。

例 9.13 如果将下推列表的长度看做是 P 所做移动次数的函数, 就能够得到如图 9-17 所示的曲线。在这幅图上, 用每次移动后 P 的格局 (而不是 ID) 来标记每个点。

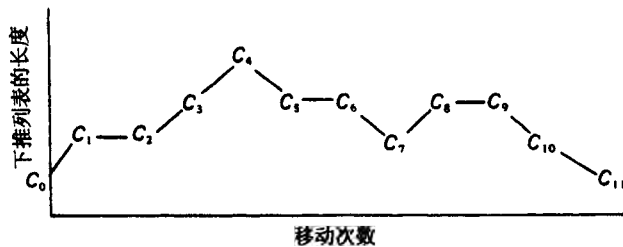


图 9-17 格局序列图

例如, 从图 9-17 中可以看到 C_0 和 C_{11} 都有 Z_0 在栈顶, 而中间格局都没有 Z_0 在栈顶, 所以可得出结论 $C_0 \Rightarrow C_{11}$ 。如果 C_{11} 是最终格局, 则无论对它自己还是对 C_0 而言, C_{11} 都是终结者。从这个图上还可以推导出 $C_1 \Rightarrow C_2$, $C_2 \Rightarrow C_7$, $C_1 \Rightarrow C_{10}$, 因为 C_{10} 造成 P 从下推列表中弹出一个符号, 所以 C_{10} 是最终格局。□

下面两个简单引理是模拟算法的关键。

引理 9.1 当且仅当存在某个 i , $0 \leq i \leq |w| + 1$, 使形如 (s_i, i, Z_0) 的格局是初始格局 $(s_0, 1, Z_0)$ 的终结者时, P 接受 w 。

证明: 这个结论可以直接从 2DPDA 如何接受输入字符串的定义得到。□

定义 如果对格局 C 和格局 D , 有 $C \Rightarrow D$, 则定义 C 是 D 的前驱 (predecessor)。如果 $C \Rightarrow D$, 则定义 C 是 D 的直接前驱 (immediate predecessor)。如果格局对 (C, D) 和格局对 (C', D') 满足下面四个条件:

$$1) C = (s, i, A) \quad D = (t, j, A)$$

⊖ 如果 $m=0$, 那么 $C \vdash C'$ 。

$$2) C' = (s', i', B) \quad D' = (t', j', B)$$

$$3) (s, i, A) \vdash (s', i', BA)$$

且

$$4) (t', j', BA) \vdash (t, j, A)$$

也就是说, 如果 P 能够通过一个入栈 (push) 移动操作从格局 C 转到 C' , 通过一个出栈 (pop) 移动从格局 D' 转到 D , 则称格局对 (C, D) 位于格局对 (C', D') 之下 (below)。(注意, (C, D) 和 (C', D') 可以不同也可以相同。)

引理 9.2 如果格局对 (C, D) 位于 (C', D') 之下, 且 $C' \Rightarrow D'$, 则 $C \Rightarrow D$ 。

证明: 证明很容易, 留作练习。 \square

例 9.14 在图 9-17 中, (C_3, C_5) 位于 (C_4, C_4) 之下, 且 (C_7, C_{10}) 位于 (C_8, C_9) 之下。然而, 因为 C_3 和 C_9 在下推列表的顶端可能没有同样的下推符号, 所以不能确定 (C_2, C_{10}) 是否位于 (C_3, C_9) 之下。 \square

模拟算法通过找到输入 w 时, 2DPDA P 的每个格局的终结者。一旦找到初始格局 $(q_0, 1, Z_0)$ 的终结者, 算法就结束了。

用数组 TERM 存储每个格局的终结者。假设格局已经 (按某种字典序) 排成线性有序的序列了。接着用把格局名 C 看成整数, 用 TERM[C] 表示 C 的终结者。

另外, 还要用到列表数组 PRED。PRED 以格局作为下标, PRED[D] 存放满足 $C \Rightarrow D$ 的所有格局 C 的列表。

除了数组 TERM 和 PRED 外, 还将使用另外两个列表 NEW 和 TEMP。列表 NEW 包含还没有被考虑的使得 TERM[C] = D 的格局对 (C, D) , 列表 TEMP 用在过程 UPDATE(C, D) 中, 用来存放格局 C 的前驱。

执行如下过程。首先查看 C 是否为最终格局。如果是, 则设置 TERM[C] = C (每个最终格局都是它自己的终结者), 且把 (C, C) 加到 NEW 中。然后考虑在 P 的一次移动中有满足 $C \vdash D$ 的格局对 (C, D) (注意: 在这样一次移动中, 下推列表不变)。

如果已知 D 的终结者, 则对于所有已知是 C 前驱的 E , 包括 C 自己, 设置 TERM[E] = TERM[D] (可以在 PRED[C] 上找到 C 的前驱)。同时, 将格局对 $(E, \text{TERM}[D])$ 添加到列表 NEW 中。

如果还不知道 D 的终结者, 那么将 C 存入 D 的直接前驱列表 PRED[D] 上。

此时, 对于每个格局 C , 无需操作栈就可以确定满足 $C \Rightarrow D$ 的唯一的格局 D , 并且 D 是 C 的终结者或者对于 $|\alpha| = 2, D \vdash (s, i, \alpha)$ 。现在考虑所有已经加到列表 NEW 的格局对。一般而言, NEW 包含尚未考虑的格局对 (A, B) , 其中, $A \Rightarrow B, B$ 是最终格局。假设 NEW 包含格局对 (A, B) , 从 NEW 中移走格局对 (A, B) , 考虑位于 (A, B) 之下的每个格局对 (C, D) 。如果 D 的终结者已经被计算, 那么设置 TERM[C] = TERM[D], 且添加格局对 $(C, \text{TERM}[D])$ 到列表 NEW 中。对于 PRED[C] 上的每个 E , 设置 TERM[E] = TERM[D], 且添加 $(E, \text{TERM}[D])$ 到列表 NEW 中。然而, 如果 D 的终结者还没有计算, 那么将 C 放到列表 PRED[D] 中。继续这个过程, 直到 NEW 为空为止。此时, 对于每个格局 C , 只要终结者存在, 就一定能确定终结者。

一旦 NEW 为空, 观察 TERM[C_0], 其中 C_0 是初始格局。如果 TERM[C_0] 是一个接受格局, 则可知 2DPDA P 接受 w , 否则, P 拒绝 w 。

下面的算法将更加精确地给出具体细节。

算法 9.4 模拟 2DPDA 算法。

输入: 一个 2DPDA $P = (S, I, T, \delta, s_0, Z_0, s_f)$ 和一个输入串 $w \in I^*, |w| = n$ 。

输出: 如果 $w \in L(P)$, 则答案是“yes”; 否则答案是“no”。

方法:

- 1) 初始化数组和列表如下: 对每个格局 C , 设 $TERM[C] = \emptyset$ 和 $PRED[C] = \emptyset$ 再设 $NEW = \emptyset$, $TEMP = \emptyset$ 。
- 2) 对于每个最终格局 C , 设 $TERM[C] = C$ 且添加格局对 (C, C) 到 NEW 中。
- 3) 对于每个格局 C , 确定在一次移动后是否有某个格局 D , 使 $C \vdash D$ 。如果有, 则调用过程 $UPDATE(C, D)$ 。过程 $UPDATE$ 如图 9-18 所示。
- 4) 当 NEW 不空时, 从 NEW 中移除格局对 (C', D') , 对于每对位于 (C', D') 之下的 (C, D) , 调用过程 $UPDATE(C, D)$ 。
- 5) 如果 $TERM[C_0]$ 是一个最终格局, 这里 C_0 是初始格局, 那么回答“yes”; 否则回答“no”。 □

例 9.15 考虑在如图 9-17 所示的 2DPDA 上运用算法 9.4 时发生的一些计算。

```

procedure UPDATE(C, D):
begin
  comment 无论何时 UPDATE(C, D) 被调用, 均有  $C \vdash D$ ;
1.  if  $TERM[D] = \emptyset$  then 将  $C$  添加到  $PRED[D]$  中
    else
      begin
2.         $TEMP \leftarrow \{C\}$ ;
3.        while  $TEMP \neq \emptyset$  do
          begin
4.            从  $TEMP$  中选择和移走格局  $B$ ;
5.             $TERM[B] \leftarrow TERM[D]$ ;
6.            把  $(B, TERM[D])$  加到  $NEW$  上;
7.            for  $PRED[B]$  中的每个  $A$  do 把  $A$  添加到  $TEMP$  中
          end
        end
      end
end

```

图 9-18 UPDATE 过程

步骤 2 确定 C_4, C_6, C_9, C_{10} 和 C_{11} 是最终格局, 因此它们是自己的终结者。我们添加格局对 $(C_4, C_4), (C_6, C_6), (C_9, C_9), (C_{10}, C_{10})$ 和 (C_{11}, C_{11}) 到 NEW 中。

步骤 3 调用过程 $UPDATE(C_1, C_2)$ 。这时, 因为 $TERM[C_2] = \emptyset$, $UPDATE$ 只把 C_1 放到列表 $PRED[C_2]$ 中。在步骤 3, 还调用 $UPDATE(C_5, C_6)$ 。因为 $TERM[C_6] = C_6$, 所以 $UPDATE$ 设 $TERM[C_5] = C_6$, 且添加 (C_5, C_6) 到 NEW 中。在步骤 3 里还调用了 $UPDATE(C_8, C_9)$, 它设 $TERM[C_8] = C_9$, 且添加 (C_8, C_9) 到 NEW 。这样, 在步骤 3 执行完后, NEW 包含如下的七个格局对:

$(C_4, C_4), (C_6, C_6), (C_9, C_9), (C_{10}, C_{10}), (C_{11}, C_{11}), (C_5, C_6), (C_8, C_9)$

在步骤 4 中, 因为 (C_3, C_3) 位于 (C_4, C_4) 之下, 所以从 NEW 中移除 (C_4, C_4) , 并调用 $UPDATE(C_3, C_3)$ 。因为此时 $TERM[C_3] = C_6$, $UPDATE$ 设 $TERM[C_3] = C_6$, 并且添加 (C_3, C_6) 到 NEW 中。然后第 4 步从 NEW 中移除 (C_6, C_6) , 因为(假设) (C_6, C_6) 之下没有格局对, 则不调用 $UPDATE$ [⊖]。类似地, 对于格局对 $(C_9, C_9), (C_{10}, C_{10}), (C_{11}, C_{11})$ 和 (C_5, C_6) , 我们均不调用 $UPDATE$ 过程。

当从 NEW 中移除 (C_8, C_9) 时, 调用 $UPDATE(C_7, C_{10})$, 使 $TERM[C_7] = C_{10}$, 并添加 (C_7, C_{10}) 到 NEW 中。此时, NEW 包含 (C_3, C_6) 和 (C_7, C_{10}) 。

当从 NEW 中移除 (C_3, C_6) 时, 调用 $UPDATE(C_2, C_7)$, 因为 $PRED[C_2]$ 包含 C_1 , 所以有 $TERM[C_2] = C_{10}$ 和 $TERM[C_1] = C_{10}$ 。把 (C_2, C_{10}) 和 (C_1, C_{10}) 添加到 NEW 中。

⊖ 为了简化问题, 假设 P 不包括图 9-17 未包含的移动。

请读者自行完成模拟过程。 □

定理 9.10 算法 9.4 在 $O(|w|)$ 时间内正确回答问题“是否 $w \in L(P)$?”

证明: 可以证明, 没有一个格局会重复(多于一次)放到 TEMP 上。这样, 每次对 UPDATE 的调用总能终止。也可以证明没有格局对被重复(多于一次)放到列表 NEW 上, 所以算法本身总会终止。作为练习, 这两部分的证明细节留给读者。

下面将证明, 在当且仅当 $w \in L(P)$ 时, 算法 9.4 的结论中, $TERM[C_0]$ 是最后格局。同时, 对算法 9.4 的执行步数进行归纳, 很容易证明以下结论:

- 1) 如果 $TERM[C]$ 设成 D , 那么 D 是 C 的终结者。
- 2) 如果将 C 添加到 $PRED[D]$, 那么 $C \Rightarrow D$ 。
- 3) 如果 (C, D) 被放到 NEW 上, 那么 $TERM[C] = D$ 。

这样, 如果算法 9.4 发现 $TERM[C_0]$ 是最后格局, 那么根据引理 9.1, $w \in L(P)$ 为真。

反之, 必须证明, 如果 D 是 C 的终结者, 那么 $TERM[C]$ 最终被设成 D 。证明可通过对序列 $C \vdash^* D$ 所包含的移动次数进行归纳得出。归纳基础(0 步移动)是显然的, 因为 $C = D$, 在算法 9.4 的第 2 步 $TERM[C]$ 设成 D 。

在归纳步骤中, 假设 $C \vdash^* E \vdash^* D$, 有两种情况要考虑。

情况 1 E 是一个格局, 即移动 $C \vdash E$ 并不是一个 push 或 pop 移动, 因此在步骤 3, $UPDATE(C, E)$ 被调用。如果此时 $TERM[E]$ 已经被设为 D , C 将在第 2 行被放到 TEMP 中, 最终 $TERM[C]$ 将在第 5 行被设为 D 。如果 $TERM[E]$ 还没有被设为 D , 那么在 $UPDATE(C, E)$ 的第 1 行, 将 C 添加到 $PRED[E]$ 中。根据归纳假设, 最终设置 $TERM[E] = D$ 。如果这个设置发生在 $UPDATE$ 的第 5 行, 那么 C 将在第 7 行被添加到 TEMP, $TERM[C]$ 在调用 $UPDATE$ 时被设置成 D 。因为 $E \neq D$, 所以 $TERM[E]$ 不可能在算法 9.4 的第 2 步被设为 D (如果 $E = D$, 那么在第 3 步考虑 $C \vdash E$ 时, $TERM[E]$ 已经被设为 D)。因此得出结论, 在第一种情况下, $TERM[C]$ 被设为 D 。

情况 2 E 是一个 ID, 且 $C \vdash E$ 是一个 push 移动, 那么能够找到格局 A, B 和 F , (C, F) 位于 (A, B) 之下, 使得 $A \vdash^* B, F \vdash^* D$, 其每个序列所含的移动都比 $C \vdash^* D$ 的序列所含移动次数少 (格局 A 是 ID E 的“表面”格局)。根据归纳假设, $TERM[A]$ 被设为 B , $TERM[F]$ 被设为 D 。

假设后者在前者之前发生, 即 $TERM[F]$ 先被设置为 D , 而后 $TERM[A]$ 被设置为 B , 那么 (A, B) 最终被放到 NEW 中, 在步骤 4, $UPDATE(C, F)$ 被调用。因为此时 $TERM[F] = D$, 所以在算法第 5 行设 $TERM[C]$ 为 D 。

反之, 假设在 $TERM[F]$ 被设为 D 之前 $TERM[A]$ 已经被设为 B , 那么当 $UPDATE(C, F)$ 被调用时, $TERM[F] = \emptyset$, 此时 C 被加到 $PRED[F]$ 。但是, 当 $TERM[F]$ 被计算时, $TERM[C]$ 被设成 D 。这样, 完成归纳并得出算法 4 的正确性证明。

考虑算法 9.4 的运行时间。因为有 $O(n)$ 个格局, 所以步骤 1 和 2 需要 $O(n)$ 时间。因为对于每个格局, 2DPDA 至多有一个可能的移动, 所以就有至多 $O(n)$ 个格局对 (C, D) , 使得 $C \vdash D$ 。这样, 在步骤 3 至多调用 $UPDATE O(n)$ 次。

当 $C' \Rightarrow D'$, 且已经发现 C' 的终结者是 D' 时, 格局对 (C', D') 被放到列表 NEW 中。因为对于每个格局终结者是唯一的 (如果有终结者的话), 那么没有格局对会多于一次放到列表 NEW 中。因此, 放到列表 NEW 的格局对总数不会超过 $O(n)$ 。对于每个放到 NEW 上的格局对, 只有有限数量的格局对位于它之下。因为如果 (C, D) 位于 (C', D') 之下, 那么 C 和 C' 的磁头位置至多相差一格, 类似, D 和 D' 的磁头位置也至多相差一个。这样, $UPDATE$ 被调用了 $O(n)$ 次。

现在考虑用在子程序 $UPDATE$ 的总时间。可以证明, 每个格局在 $PRED$ 数组里至多出现一次, 也没有格局会多于一次被放置在列表 TEMP 中, 这样, $UPDATE$ 的第 4~6 行的执行时间可

以归结为从 TEMP 移走的格局 B 的数量, 第 7 行的执行时间可以归结为添加到 TEMP 的格局 A 的数量。因为 UPDATE 至多调用 $O(n)$ 次, 可以得出, 除了上面归结到格局 A 和 B 的时间外, UPDATE 花费的总时间是 $O(n)$ 。因此, 算法 9.4 的运行时间是线性的。□

本节所得到的结论的基本应用是证明某些问题存在线性时间的算法。我们已经看到, 一些模式匹配问题能够形式化地转换为语言识别问题。如果能设计一个 2DPDA 以接受对应于模式匹配问题的语言, 那么就on知道对该问题存在一个线性时间的算法。2DPDA 对长度为 n 的输入可能做 n^2 次甚至 2^n 次移动, 这个事实经常使在 2DPDA 上设计一个算法来识别语言(这个语言是相应于模式匹配问题的语言)比直接在 RAM 上设计一个解决模式匹配问题的线性算法更为容易。

9.5 位置树和子串标识符

上节已经证明, 如果模式匹配问题能够形式化为寻找 2DPDA 的语言识别问题, 那么就能够在线性时间内解决最初的模式匹配问题。可以直接使用算法 9.4 作为解决初始问题的一个线性时间算法。然而, 由模拟产生的常数因子使这个方法并不是非常吸引人。本节我们将讨论一种数据结构, 这种数据结构能够更加实际地应用于线性时间的模式匹配算法中。一些模式匹配问题可以应用这种数据结构, 包括:

- 1) 给定文本串 x 和模式串 y , 确定 y 在 x 中的所有出现。
- 2) 给定文本串 x , 确定 x 的最长的重复的子串。
- 3) 给定字符串 x 和 y , 确定 x 和 y 的最长公共子串。

定义 一个长度为 n 的字符串中的位置(position)是 $1 \sim n$ 之间的整数。如果 $x = yaz$, 且 $|y| = i-1$, 则称符号 a 出现在字符串 x 的位置 i 。如果 $x = yuz$, $|y| = i-1$, 除非 $y' = y$, 否则 x 不能写成 $y'uz'$, 则称子串 u 标识字符串 x 中的位置 i 。也就是说, u 在 x 中仅有的出现在位置 i 开始。例如, 子串 bba 标识字符串 $abbabb$ 的位置 2。子串 bb 则不能标识位置 2。

在本章的余下部分中, 设 $x = a_1a_2 \cdots a_n$ 是某个字母表 I 上的一个字符串, $a_{n+1} = \$$ 是一个不在 I 中的符号, 那么 $x\$ = a_1a_2 \cdots a_na_{n+1}$ 中的每个位置 i 至少被一个字符串标识, 即 $a_ia_{i+1} \cdots a_{n+1}$ 。我们把标识 $x\$$ 中位置 i 的最短字符串称为 x 中的位置 i 的子串标识符, 记作 $s(i)$ 。

例 9.16 考虑字符串 $abbabb\$$, 图 9-19 列出了其位置 1~7 的子串标识符的列表。□

对于一个字符串 $x\$$ 的各位置子串标识符, 可以方便地用称做 I 树的一棵树来表示, 树中的边和某些顶点均带有标记。

定义 一棵 I 树是一棵带标记的树 T , 对于 T 的每个内部顶点 v , 从 v 出发的边用字母表 I 中的字符做标记。如果 T 中的边 (v, w) 标记为 a , 那么称 w 为 v 的 a 孩子。

字符串 $x\$ = a_1a_2 \cdots a_na_{n+1}$ 的位置树(position tree)是一棵 $(I \cup \{\$, \})$ 树, 其中 a_i 在 I 中, $1 \leq i \leq n$, $a_{n+1} = \$$, 该树满足以下条件:

- 1) T 有 $n+1$ 个标记为 $1, 2, \cdots, n+1$ 的叶子。 T 的叶子是和 $x\$$ 中的位置一一对应的。
- 2) 沿着从根到标记为 i 的叶子的路径, 所经过的边的标记组成的序列是 $s(i)$, 即位置 i 的子串标识符。

例 9.17 字符串 $abbabb\$$ 的位置树如图 9-20 所示。例如, 从根到标记为 2 的叶子的路径拼出的字符串是 bba , 它是位置 2 的子串标识符。□

下面的引理阐述子串标识符的几个基本属性。

引理 9.3 设 $s(i)$ 是字符串 $x\$ = a_1a_2 \cdots a_na_{n+1}$ 的位置 i 的子串标识符。

- 1) 如果 $s(i)$ 长度为 j , 则 $s(i-1)$ 的长度至多为 $j+1$ 。

位置	子串标识符
1	<i>abba</i>
2	<i>bba</i>
3	<i>ba</i>
4	<i>abb\$</i>
5	<i>bb\$</i>
6	<i>b\$</i>
7	<i>\$</i>

图 9-19 $abbabb\$$ 的子串标识符

2) 没有子串标识符正好是另一个子串标识符的前缀。

证明:

a) 如果 $s(i-1)$ 的长度大于 $j+1$, 则有某个位置 $k \neq i-1$ 使 $a_{i-1}a_i \cdots a_{i+j-1} = a_k a_{k+1} \cdots a_{k+j}$ 。因此, $a_i a_{i+1} \cdots a_{i+j-1} = a_{k+1} a_{k+2} \cdots a_{k+j}$, 且 $a_i \cdots a_{i+j-1}$ 不标识位置 i , 与前提 $a_i \cdots a_{i+j-1}$ 能标识位置 i 矛盾。

b) 证明很容易, 作为练习留给读者。 □

作为引理 9.3(b) 的一个结论, 可以断定事实上每个字符串 $x\$$ 都存在一棵位置树。

许多模式匹配问题能够用位置树解决, 包括本节开始时提到的那些问题。

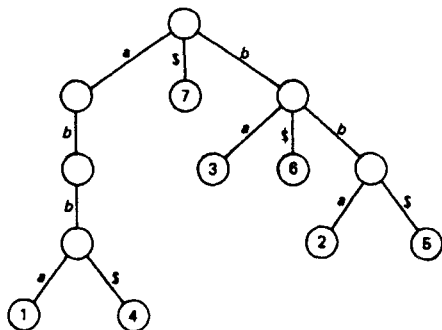


图 9-20 一棵位置树

例 9.18 考虑基本的模式匹配问题: “ $y = b_1 b_2 \cdots b_p$ 是 $x = a_1 a_2 \cdots a_n$ 的子串吗?” 假设构造了 $x\$$ 的位置树 T , 为了确定 $y = b_1 b_2 \cdots b_p$ 是否为 x 的子串, 把位置树看做一个确定型有穷自动机的转换图。也就是说, 从 T 的根开始, 在位置树中沿着最长的可能路径, 拼出 $b_1 b_2 \cdots b_j (0 \leq j \leq p)$ 。假设这个路径在顶点 v 终止, 可能发生以下三种情况。

1) 如果 $j < p$ 且顶点 v 不是叶子, 则回答 “no”。在这种情况下, 我们知道 $b_1 b_2 \cdots b_j$ 是 x 的子串, 但 $b_1 b_2 \cdots b_{j+1}$ 不是。

2) 如果 $j \leq p$ 且顶点 v 是一个标记为 i 的叶子, 那么知道 $b_1 b_2 \cdots b_j$ 匹配 x 中从位置 i 开始的 j 个字符。必须比较 $a_{i+j} a_{i+j+1} \cdots a_{i+p-1}$ 和 $b_{j+1} b_{j+2} \cdots b_p$ 。如果它们不匹配, 那么答案为 “no”; 否则回答 “yes” 且报告 y 是 x 在位置 i 开始的子串。

3) 如果 $j = p$ 且顶点 v 不是叶子, 那么回答 “yes”。在这种情况下, y 是 x 中在两个或更多位置上开始的子串, 这些位置由位置树中顶点 v 的子树所包含的叶子上的标记给出。

例 9.19 位置树能够用来在给定的串 $x = a_1 a_2 \cdots a_n$ 中找到最长的重复子串 (子串的两次出现可能有重叠)。最长的重复子串对应于最大深度的位置树的内部顶点。这样的顶点能够用一种直接的方式进行定位。

考虑在 $abbabb$ 中查找最长重复子串。在如图 9-20 的 $abbabb\$$ 的位置树上, 有一个深度为 3 的内部顶点, 且没有其他更深的内部顶点。这样相应于这个顶点的字符串 abb 是 $abbabb$ 中最长的重复子串。标记为 1 和 4 的叶子表明 abb 两次出现的开始位置, 且它们之间没有重叠。 □

现在详细考虑位置树的构造问题。本章余下部分使用 $a_1 a_2 \cdots a_n a_{n+1}$ 表示 $x\$$, 其中 a_{n+1} 是唯一的右端标识符 $\$$ 。用 x_i 表示后缀 $a_i \cdots a_n a_{n+1}$, $1 \leq i \leq n+1$; 用 $s_i(j)$ 表示 x_i 中位置 j 的子串标识符。所有的位置都是按照初始的字符串 $a_1 a_2 \cdots a_n a_{n+1}$ 中的位置顺序编号。

下面将给出构造 $a_1 a_2 \cdots a_n a_{n+1}$ 的位置树的算法, 其所需时间和结果的位置树中顶点个数成比

⊙ a^n 表示 n 个 a 连接而成的字符串。

例。注意 $a_1 a_2 \cdots a_n a_{n+1}$ 的位置树可能有 $O(n^2)$ 个顶点, 例如, $a^n b^n a^n b^n \$^\circ$ 的位置树有 $n^2 + 6n + 2$ 个顶点(对此读者可以验证)。然而, 在对“任意”字符串的含义进行合理假设的前提下(如从固定的字母表里均匀且独立地选出的字符构成的符号串), 我们能够证明长度为 n 的字符串的一棵“平均”位置树有 $O(n)$ 个顶点。尽管这里不给出证明, 但是确有算法可从一个给定的字符串直接构造出一棵压缩形式的位置树, 其所用时间最坏情况下为 $O(n)$ 。参考文献包含了对这个算法的引用。

因为需要在下面的算法中从 S_{i+1} 构造 S_i , S_i 是 x_i 子串标识符的集合, S_{i+1} 是 x_{i+1} 子串标识符的集合, 所以, 先考虑 S_i 和 S_{i+1} 之间的区别。一个明显的区别是 S_i 包含 $s_i(i)$, 即包含 x_i 的第一个位置的子串标识符。由于 S 包含这个串, 所以 S_{i+1} 的 k 位置的子串标识符可能不再是 S_i 中 k 位置的子串标识符。当且仅当 $s_{i+1}(k)$ 是 $s_i(i)$ 的前缀($s_{i+1}(k)$ 是 S_{i+1} 中 k 位置的子串标识符), 这种情况出现。在这种情况下, 必须加长 $s_{i+1}(k)$ 得到 $s_i(k)$ ($s_i(k)$ 是 S_i 中位置 k 的子串标识符)。 S_{i+1} 中不可能有两个子串标识符都需要加长。假设 $s_{i+1}(k_1)$ 和 $s_{i+1}(k_2)$ 都需要加长, 则这两个字符串将都是 $s_i(i)$ 的前缀, 因此一个将是另一个的前缀, 与引理 9.3(b) 矛盾。

位置 p	$s_4(p)$	$s_3(p)$	$s_2(p)$	$s_1(p)$
7	\$	\$	\$	\$
6	b\$	b\$	b\$	b\$
5	bb	bb	bb\$	bb\$
4	a	a	a	abb\$
3	-	ba	ba	ba
2	-	-	bba	bba
1	-	-	-	abba

图 9-21 $abbabb \$$ 的部分后缀的子串标识符

例 9.20 设 $a_1 a_2 \cdots a_n a_{n+1} = abbabb \$$ 。 $x_4 = abb \$$, $x_3 = babb \$$, $x_2 = bbabb \$$ 和 $x_1 = abbabb \$$ 的子串标识符如图 9-21 中的列表所示。 S_3 是在 S_4 上只加了一个 $S_3(3) = ba$ 得到。另一方面, S_2 需要作两个改变得到 S_3 : 增加 $s_2(2) = bba$ 到 S_2 中, 在 $s_2(5)$ 的末尾添加 \$ 得到 $s_2(5) = bb \$$ 。为了得到 S_1 , 需要增加 $s_1(1) = abba$ 到 S_2 , 在 $s_2(4)$ 上追加 $bb \$$ 得到 $s_1(4) = abb \$$ 。 □

考虑从 T_{i+1} (表示 S_{i+1} 的位置树) 构造 T_i (表示 S_i 的位置树) 涉及的相关操作。给定 T_{i+1} , 必须添加一片叶子到 T_{i+1} , 这片叶子相应于 $s_i(i)$, 标记为 i 。对 $(i < k \leq n)$, 如果 $s_{i+1}(k)$ 是 $s_i(i)$ 的前缀, 那么必须在 T_{i+1} 中延长从根到标记为 k 的叶子的路径, 所以 T_i 中标记为 k 的新叶子将相应于 $s_i(k)$ 。由此可见, 有效地从 T_{i+1} 构造 T_i (x_i 的位置树) 的关键是能够很快找到字符串 y , 使得 $a_i y$ 是 x_i 的最长前缀, 也是 x_{i+1} 的子串, 且确定 $a_i y$ 是否 T_{i+1} 中的一个子串标识符的前缀。

构造需要关联三个新的结构到一棵位置树上。第一个结构是一个位向量(数组), 将它和位置树 T_i 的每个顶点相关联, 用 B_v 表示顶点 v 的位向量, I 中的每个符号在 B_v 中有一个向量分支。在 T_i 中, 若顶点 v 对应字符串 y 且 $a \in I$, 那么, 如果 ay 是 x_i 的子串, 则 $B_v[a]$ 是 1; 否则 $B_v[a]$ 是 0。

接下来, 将位置树上的每个顶点的深度关联到顶点上。随着树逐渐长大, 这个信息很容易更新, 因此后面假设它会随之计算更新而不再具体提及这个过程。

最后需要增加到位置树的是一个树结构, 该树结构建立在位置树的顶点上。我们称这个树结构为“辅助树”, 事实上它是定义在位置树顶点上的另一个边的集合。

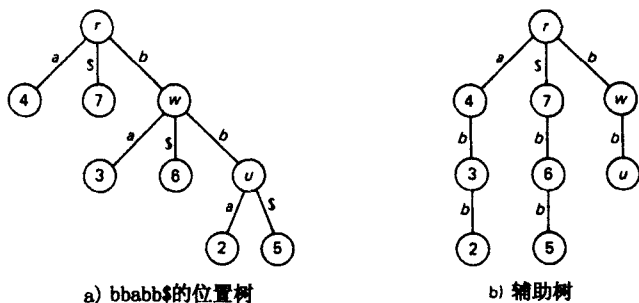


图 9-22

定义 设 T_i 是 $x_i = a_i a_{i+1} \cdots a_n \$$ 的位置树, T_i 的辅助树 (auxiliary tree) 树是一棵 $(I \cup \{ \$ \})$ 树, 满足以下条件:

1. A_i 和 T_i 的顶点集一样。
2. 如果 ay 是 T_i 中对应顶点 w 的字符串, y 是 T_i 中对应顶点 v 的字符串, 那么在 A_i 中, 顶点 w 是顶点 v 的 a 儿子, 这样在 A_i 中, 从根到 w 的路径可拼出 $y^R a$ (R 表示反转。——译者注), 而在 T_i 中从根到 w 的路径可拼出 ay 。

例 9.21 $x_2 = bbabb \$$ 的位置树和辅助树如图 9-22 所示 (叶子上的编号和 $x = abbabb \$$ 的字符位置对应)。辅助树的叶子不一定是位置树的叶子, 但是它们的顶点集合是一样的。□

从定义不能很明显地看出一棵位置树是否有辅助树。事实上, 对于一棵任意的 I 树, 可能没有辅助树。下面的定理给出了一棵树拥有辅助树的条件。

定理 9.11 一棵 I 树 T 拥有辅助树的充分必要条件是, 如果 T 中有一个顶点相应于字符串 ax , 这里 $a \in I, x \in I^*$, 那么也有一个顶点相应于字符串 x 。

证明: 留作练习。□

推论 每棵位置树都有辅助树。

证明: 如果 ax 是位置 i 的子串标识符的前缀, 那么根据引理 9.3(a), x 是位置 $i+1$ 的子串标识符的前缀。□

现在准备给出从 T_{i+1} (x_{i+1} 的位置树) 计算 T_i (x_i 的位置树) 的算法。

算法 9.5 从 T_{i+1} 构造 T_i 。

输入: 一个串 $a_1 \cdots a_n a_{n+1}$, 字符串 $x_{i+1} = a_{i+1} \cdots a_n a_{n+1}$ 的位置树 T_{i+1} 和辅助树 A_{i+1} 。

输出: 字符串 x_i 的位置树 T_i 和辅助树 A_i 。

方法:

1. 在 T_{i+1} 中找到标记为 $i+1$ 的叶子 (这是最后加到 T_{i+1} 的叶子);
2. 遍历从这片叶子到 T_{i+1} 的根的路径直到遇到顶点 u_y , 使得 $B_{u_y}[a_i] = 1$ (顶点 u_y 相应于最长的串 y , 使得 $a_i y$ 是 x_i 的一个前缀, 也是 x_{i+1} 的位置 k 开始的一个子串, $i < k \leq n$)。如果没有这样的顶点存在, 继续遍历直到根;
3. 在从标记为 $i+1$ 的叶子到 u_y 的儿子的路径上的每个顶点 v , 设 $B_v[a_i] = 1$; 如果 u_y 不存在, 则对从标记 $i+1$ 的叶子到根的路径上的每个顶点 v , 设置 $B_v[a_i] = 1$ 。(在这条路径上的每个顶点 v 相应于 x_{i+1} 的某个前缀 z 。因此, 很明显 $a_i z$ 是 $a_i x_{i+1}$ 的子串。)
4. i) 如果 u_y 不存在, 则转向情况 1;
 ii) 否则, 如果在辅助树 A_{i+1} 上 u_y 没有 a_i 儿子, 则转向情况 2;
 iii) 否则, 如果在辅助树 A_{i+1} 中, u_y 有一个 a_i 儿子 v_y , 则转向情况 3。

顶点 v_i 相应于位置树 T_{i+1} 中的 $a_i y$ 。

情况 1 a_i 是不出现在 x_{i+1} 中字符, 这样, 位置 i 的子串标识符就是 a_i 。为了从 T_{i+1} 构造 T_i , 作如下操作。

1) 创建一个新叶子, 标记为 i , 作为 T_{i+1} 的根的 a_i 儿子;

2) 对所有 $a \in I$, 设 $B_i[a] = 0$ 。

为了从 A_{i+1} 构造 A_i , 将标记为 i 的新顶点作为 A_{i+1} 的根的 a_i 儿子。

情况 2 a_i 出现在 x_{i+1} 中, 但是只有 $a_i y$ 的某个前缀出现在 T_{i+1} 中 (拼出一条从 T_{i+1} 的根到某片叶子 v_1 的路径)。这种情况出现在 T_{i+1} 中位置 k ($i < k \leq n$) 的子串标识符必须延长, 才能成为 T_i 中位置 k 的子串标识符时。假设 $y = a_{i+1} a_{i+2} \cdots a_j$, 叶子 v_1 相应于 $a_i a_{i+1} \cdots a_p$, $p < j$, 那么, 叶子 v_1 标记为 k , $a_i a_{i+1} \cdots a_p$ 是 T_{i+1} 中位置 k 的子串标识符 (即 $a_i a_{i+1} \cdots a_p = a_k a_{k+1} \cdots a_l$)。

为了从 T_{i+1} 构造 T_i , 增加一棵子树到顶点 v_1 上, 这棵子树带着两片新叶子, 分别标记为 i 和 k 。在这棵增加的子树中, 从 v_1 到 k 的路径拼出字符串 $a_{i+1} a_{i+2} \cdots a_{m+1}$, 从 v_1 到 i 的路径拼出字符串 $a_{p+1} a_{p+2} \cdots a_{j+1}$, 这里 $a_{i+1} a_{i+2} \cdots a_m = a_{p+1} a_{p+2} \cdots a_j$ 。这样, 从 T_i 的根到叶子 k 的路径将拼出 $a_k a_{k+1} \cdots a_m a_{m+1}$, 这是 x_i 中位置 k 的新子串标识符, 从根到叶子 i 的路径将拼出 $a_i a_{i+1} \cdots a_j a_{j+1}$, 这是 x_i 中位置 i 的子串标识符。

为了从 T_{i+1} 构造 T_i , 具体步骤如下。

1) 在 T_{i+1} 中沿着从 u_y 到根的路径找到 u_1 , u_1 是 A_{i+1} 中有 a_i 儿子的 u_y 的第一个祖先, 记这个 a_i 儿子为 v_1 。在 T_{i+1} 中, 顶点 v_1 是标记为 k 的叶子, $i < k \leq n$ 。从顶点 v_1 中移走标记 k ;

2) 设 $u_1, u_2, \cdots, u_q, u_y$ 是 T_{i+1} 中从 u_1 到 u_y 路径上的顶点序列。假设从 u_h 到 u_{h+1} 的分支标记为 c_h , $1 \leq h < q$, 从 u_q 到 u_y 的分支标记为 c_q , 如图 9-23 所示 (在上面讨论中, $c_1 c_2 \cdots c_q = a_{p+1} a_{p+2} \cdots a_j = a_{i+1} a_{i+2} \cdots a_m$);

3) 在 T_i 中, 创建 q 个新顶点 $v_2, v_3, \cdots, v_q, v_y$, 使 v_{h+1} 成为 v_h 的 c_h 儿子 ($1 \leq h < q$), 并使 v_y 是 v_q 的 c_q 儿子;

4) 设 j 为 i 与 u_y 的深度之和, m 是 k 与 u_y 的深度之和, 给 T_i 两个标记为 i 和 k 的新叶子, 使叶子 i 是 v_y 的 a_{j+1} 儿子, 叶子 k 是 v_y 的 a_{m+1} 儿子;

5) 对所有的 $a \in I$ 和所有 $\{v_2, v_3, \cdots, v_q, v_y, k\}$ 中的 v , 令 $B_v[a] = B_{u_1}[a]$;

6) 对所有 $a \in I$, 令 $B_i[a] = 0$ 。

为了从 A_{i+1} 构造 A_i , 需作如下工作。

7) 使 u_y 成为 A_i 中的 v_y 的 a_i 儿子;

8) 设 u' 是 T_{i+1} 中 u_y 的 a_{j+1} 儿子, 使标记为 i 的新叶子成为 A_i 中 u' 的 a_i 儿子;

9) 设 u'' 是 T_{i+1} 中 u_y 的 a_{m+1} 儿子, 使标记为 k 的新叶子成为 A_i 中 u'' 的 a_i 儿子;

10) 使 v_h 成为 A_i 中 u_h 的 a_i 儿子, $2 \leq h \leq q$ 。

情况 3 在 A_{i+1} 中 u_y 有一个 a_i 儿子 v_y , 分两种情况考虑。

a) v_y 是 T_{i+1} 中标记为 k 的叶子, 这种情况发生在 $s_i(i) = a_i a_{i+1} \cdots a_j a_{j+1}$ 和 $s_{i+1}(k) = a_k a_{k+1} \cdots a_m$ 时, 这里 $a_k a_{k+1} \cdots a_m = a_i a_{i+1} \cdots a_j$, 这是情况 2 的特例, 此时, $v_1 = v_y$ 。为了从 T_{i+1} 构造 T_i , 从 v_y 移走标记 k , 给 v_y 一个标记为 i 的 a_{j+1} 儿子, 和标记为 k 的 a_{m+1} 儿子。该构造细节和情况 2 中除去步骤 3、7 和 10 后是一样的。

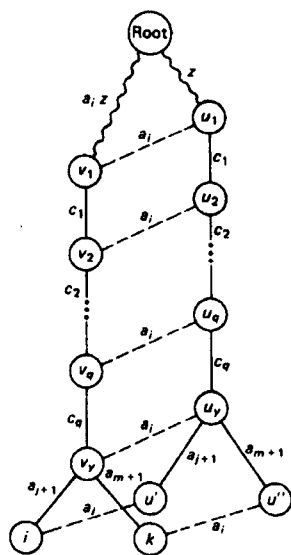


图 9-23 情况 2 中 T_i 的重要部分, 虚线表示 A_i 中的边

b) v_y 是 T_{i+1} 中的内部顶点, 这种情况发生在 $S_i = S_{i+1} \cup \{s_i(i)\}$ 时。为了从 T_{i+1} 构造 T_i , 可简单地给 v_y 一个 a_{j+1} 儿子 (v_y 原先没有 a_{j+1} 儿子), 且标记这个叶子为 i 。详细的算法如下。

1) 设 j 是 i 与 u_y 的深度的和;

2) 在 T_i 中, 使标记为 i 的新顶点成为 v_y 的 a_{j+1} 儿子 (注意, 由 y 的极大性, 即 y 是使得 $a_i y$ 是 x_i 的最长前缀, 在 T_{i+1} 中 v_y 不可能有 a_{j+1} 儿子);

3) 对所有的 $a \in I$, 设 $B_i[a] = 0$ 。因为 $a_i y a_{j+1}$ 不是 x_{i+1} 的子串, 因此, 对任意 a , $aa_i y a_{j+1}$ 一定不是 x_i 的子串;

4) 为了从 A_{i+1} 构造 A_i , 在 T_{i+1} 中使顶点 u' 成为 u_y 的 a_{j+1} 儿子 (在 T_{i+1} 中, u' 相应于 ya_{j+1}), 使新顶点 i 成为 A_{i+1} 中 u' 的 a_i 儿子 (在 A_i 中, i 将对应于 $a_{j+1} y^R a_i$)。

在情况 3(b) 中提到的顶点之间的关系如图 9-24 所示。

□

例 9.22 设 $a_1 \cdots a_n a_{n+1} = abbabb \$$, 设 T_2 和 A_2 是图 9-22 中的树。从 T_2 和 A_2 , 应用算法 9.5 构造 T_1 和 A_1 。因此, 这里 $i = 1$, $a_1 = a$ 。

首先, 定位标记为 2 的叶子, 因为 $B_2[a] = 0$, 所以设 $B_2[a] = 1$, 并且在图 9-22 中向上移动到 u 顶点, 在顶点 u , 我们发现 abb 是 $bbabb \$$ 的子串, 所以 $B_u[a] = 1$ 。这样, 顶点 u 是 u_y 。现在找出 v_y , 即 A_2 中 u 的 a 儿子。由于 v_y 不存在, 所以, 属于情况 2。

T_2 中 u 的父顶点 w 在 A_2 中没有 a 儿子。然而, w 的父顶点 r 有顶点 4 作为其 a 儿子。因此, 在情况 2 的第 1 步, 我们发现 v_1 是 T_2 和 A_2 中标记为 4 的顶点。在步骤 2, 我们发现 $q = 2$, $u_1 = r$ 和 $u_2 = w$, 另外, $c_1 = c_2 = b$, 因此在步骤 3, 创建一个新点 v , 使 v 是 T_2 中标记为 4 的顶点的 b 儿子 (在 T_1 中顶点 4 已经不使用这个标记, 而用 v_1)。同时也创建了 v_y , 并使其成为 v 的 b 儿子。

在第 4 步, 计算得到 $j = 3$ 。因为顶点 v_1 以前的标记 k 是 4, 所以算得 $m = 6$ 。我们发现 $a_{j+1} = a$, $a_{m+1} = \$$, 这样, 标记为 1 和 4 的新顶点分别作为 v_y 的 a 儿子和 $\$$ 儿子。树 T_1 如图 9-25 所示。余下的步骤留作练习。

□

引理 9.4 如果 T_{i+1} 和 A_{i+1} 是 x_{i+1} 的位置树和辅助树, 那么, 根据算法 9.5 构造的 T_i 和 A_i 是 x_i 的正确的位置树和辅助树。

证明: 假设 T_{i+1} 是 S_{i+1} 的表示, S_{i+1} 是 x_{i+1} 中各位置的子串标识符的集合, A_{i+1} 是 T_{i+1} 的辅助树, 则 S_i 有两种可能:

1. $S_i = S_{i+1} \cup \{s_i(i)\}$;
2. $S_i = S_{i+1} - \{s_{i+1}(k)\} \cup \{s_i(k), s_i(i)\} (i < k \leq n)$ 。

第一种可能性由算法 9.5 的情况 1 和情况 3(b) 所覆盖。第二种可能性包含在情况 2 和 3(a) 中, 情况 1 ~ 3 中需要的算法正确性的证明和推导包含在算法的描述中。

□

可使用如下过程构造任意字符串的位置树。

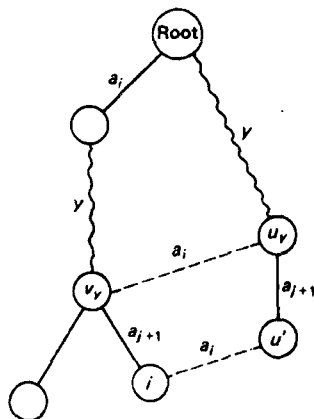


图 9-24 情况 3(b) 中 T_i 的重要部分, 两条标记 a_i 的虚线是 A_i 中的边

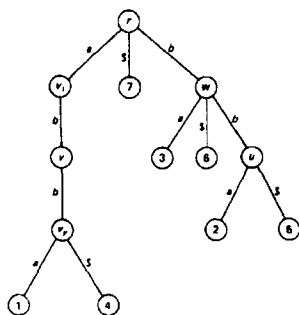
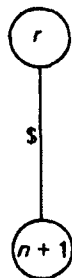
图 9-25 最后的位置树 T_1 

图 9-26 位置树的初始图

算法 9.6 构造一棵位置树。

输入：一个字符串 $x\$ = a_1 \cdots a_n a_{n+1}$, $a_i \in I$, $1 \leq i \leq n$ 且 $a_{n+1} = \$$ 。

输出： $x\$$ 的位置树 T_1 。

方法：

1. 令 T_{n+1} 是如图 9-26 所示的位置树，且对所有的 $a \in I$ ，令 $B_i[a] = B_{n+1}[a] = 0$ ；
2. 令 A_{n+1} 是和图 9-26 一样的辅助树；
3. for $i \leftarrow n$ step-1 until 1 do 用算法 9.5 从 T_{i+1} 和 A_{i+1} 构造 T_i 和 A_i 。

□

定理 9.12 在与位置树 T_1 中的顶点数成比例的时间内，算法 9.6 构造 $x\$$ 的位置树。

证明：算法 9.5 的步骤 1 能够在固定时间内找到叶子 $i+1$ ，将叶子 $i+1$ 加到 T_{i+1} 时保持指向该叶子的指针，当把叶子 i 加到 T_i 中时，设这个指针指向叶子 i 。

步骤 2 和 3 每次执行的工作量明显地与从顶点 $i+1$ 遍历到 u_i 的路径长度成比例，步骤 1 的每次执行的工作量是常量。情况 1~3 中任何一种情况的执行需要的时间都和添加到树中的顶点数成比例，对此，读者可以容易地进行验证。这样，算法 9.5 的各部分的执行时间除了步骤 2 和 3 外都是和 T_1 的规模成比例的。

我们仍需证明在 T_{i+1} 中 ($1 \leq i \leq n$)，顶点 $i+1$ 和 u_i (如果 u_i 不存在，则是顶点 $i+1$ 和根) 之间的距离的和不大过 T_1 的规模。记这些距离为 d_2, d_3, \dots, d_{n+1} ，用 e_i ($1 \leq i \leq n+1$) 表示 T_1 中顶点 i 的深度，对情况 1~3 作简单的分析可得：

$$e_i \leq e_{i+1} - d_{i+1} + 2 \quad (9-1)$$

如果计算 i 从 1 取到 n 的式 (9-1) 两边的和，可发现

$$\sum_{i=2}^{n+1} d_i \leq 2n + e_{n+1} - e_1 \quad (9-2)$$

从式 (9-2) 立即得到：算法 9.5 的步骤 2 和 3 的执行时间是 $O(n)$ ，因此，该时间至多和 T_1 的规模成比例。 □

正如前面提到的，一个长度为 n 的字符串的位置树可能有 $O(n^2)$ 个顶点，因此，任意模式匹配算法构造这棵树的时间复杂度都是 $O(n^2)$ 。然而，可以通过将位置树里所有链浓缩成单个顶点来“压缩”位置树和辅助树。一条链是一条路径，其所有顶点恰好有一个儿子。不难证明，对于长度为 n 的字符串，其压缩位置树至多有 $4n-2$ 个顶点。压缩树能够表示和初始的位置树一样的信息，这样，能够用在同样的模式匹配算法中。

对算法 9.5 进行修改，能够在线性时间内产生压缩的位置树和辅助树。这里没有给出修改算法，是因为它和算法 9.5 很类似，事实上在许多应用中，会希望位置树的规模和输入串的长度线性成比例。在参考文献中给出了关于位置树的压缩算法及其一些应用的资料。

有一个是 x 的子串};

b) $\{xyy^R \mid x \text{ 和 } y \text{ 是 } I^*, |y| \geq 1\}$;

c) $\{x_1cx_2\cdots cx_n \mid x_i \in \{a, b\}^*, 1 \leq i \leq n, \text{ 所有的 } x_i \text{ 是不同的}\}$;

d) $\{x_1cx_2\cdots cx_kdy_1cy_2\cdots cy_l \mid \text{每个 } x_i \text{ 和 } y_i \text{ 在 } \{a, b\}^* \text{ 中, 且对某个 } i \text{ 和 } j, \text{ 有 } x_i = y_j^R.\}$

9.16 考虑满足如下规则的 2DPDA P :

$\delta(s_0, a, Z_0) = \delta(s_0, a, A) = (s_0, +1, \text{push } A),$

$\delta(s_0, \$, A) = (s_1, -1),$

$\delta(s_1, a, A) = \delta(s_1, a, Z_0) = (s_1, -1, \text{pop}),$

$\delta(s_1, \epsilon, Z_0) = (s_f, 0).$

a) 试列出输入 aa 的 P 的所有表面格局,

b) 试使用算法 9.4 计算数组 TERM。

*9.17 假设 2DPDA 在一个移动中能够从位置 i (n 是输入长度) 移到如下的每个位置。

a) 到位置 $n-i$

b) 到位置 $i/2$

c) 到位置 $n/2$

d) 到位置 $\log n$

试证明这些修改中没有一个是增加了 2DPDA 的识别能力。考虑能否在不增加 2DPDA 的接受能力的情况下再增加一些其他能力。

**9.18 试构造一个 2DPDA 接受语言

$$\{w_1w_1^Rw_2w_2^Rw_3w_3^Ru \mid w_1, w_2, w_3 \in \{a, b\}^*, u \in \{a, b\}^*\}.$$

[提示: 设 x 是输入串, 写一个子程序找到最短的 w_1 , 使 $w_1w_1^Ry = x$ 。然后递归运用这个子程序以确定 y 左部的符号。]

**9.19 试给出识别如下语言的线性时间算法:

a) $\{ww^R \mid w \in I^*\}^*$;

b) $\{ww^Ru \mid w, u \in I^*, |w| \geq 1\}$;

c) $\{w_1w_2\cdots w_n \mid w_i \in \{a, b\}^*, \text{ 每个 } w_i \text{ 是一个非空的偶数长度的回文}\}.$

9.20 试构造一个线性算法对字符串 $a_1a_2\cdots a_n$ 和 $b_1b_2\cdots b_p$, 找到最长的 l , 使 $a_1a_2\cdots a_l$ 是 $b_1b_2\cdots b_p$ 的子串。并考虑算法怎样用来测试一个给定的字符串是否回文。

下面四个习题涉及下推自动机的一些变形。

*9.21 k 磁头 2DPDA 是在输入带上有 k 个独立的读磁头的 2DPDA。试证明能够在 $O(n^k)$ 时间内确定长为 n 的输入字符串可否被 k 磁头 2DPDA 所接受。

9.22 单向下推自动机从不向左移动输入磁头。非确定型下推自动机在每步移动时都有 0 个或多个选择。这样, 下推自动机有四种类型: 1DPDA, 1NPDA, 2DPDA 和 2NPDA。1NPDA 接受的语言叫做上下文无关语言。

a) 试证明 $\{wuw^R \mid w \in \{a, b\}^*\}$ 能够被 1DPDA 所接受;

b) 试证明 $\{ww^R \mid w \in \{a, b\}^*\}$ 能够被 1NPDA 所接受; (这个语言不可能被任何 1DPDA 所接受。)

c) 试证明 $\{wvx \mid w, x \in \{a, b\}^*, |w| \geq 1\}$ 能够被 2NPDA 所接受。(这个语言不可能被任何 1NPDA 所接受。)

*9.23 试证明语言 L 由乔姆斯基上下文无关语法范式所产生, 当且仅当 L 可被一个 1NPDA 所接受。

*9.24 试证明能够在 $O(n^3)$ 时间内确定长度为 n 的输入串能否被一个 2NPDA 所接受。

- 9.25 试为以下的字符串找到位置树。
 a) $baaaab \$$;
 b) $abababa \$$ 。
- 9.26 试证明 $a^n b^n a^n b^n \$$ 的位置树有 $n^2 + 6n + 2$ 个顶点。
- *9.27 试证明, 如果所有位置的字符都是从固定的字母表中独立且均匀地挑选出来的, 那么随机输入字符串 x 的位置树有 $O(|x|)$ 个顶点。
- 9.28 试为习题 9.25 中的位置树找到:
 a) 辅助树;
 b) 每个顶点 v 的位向量 B_v 。
- 9.29 试证明每棵位置树有一棵辅助树。
- 9.30 试通过证明能够从 T_{i+1} 和 A_{i+1} 正确地构造 T_i 和 A_i 来完成引理 9.4 的证明。
- *9.31 试说明 x 的位置树是怎样用来测试 y_1, y_2, \dots, y_m 中的任何一个是否 x 的子串, 且所需时间与 $|y_1| + |y_2| + \dots + |y_m|$ 成比例。
- 9.32 设计一个算法, 找到两个已知字符串 x 和 y 的最长公共子串, 并分析算法的时间复杂度。
- 9.33 对于给定的字符串 x 和 $b_1 b_2 \dots b_p$, 设计一个有效的算法为每个 $i, 1 \leq i \leq p$, 找到是 $b_i b_{i+1} \dots b_p$ 的前缀的 x 的最长子串。
- 9.34 给定字母表 I 上的字符串 $x = a_1 a_2 \dots a_n$ 和 $y = b_1 b_2 \dots b_p$, 设计一个有效的算法找到 y 的形如字符串 $c_1 c_2 \dots c_m$ 的最短表示, 这里每个 c_i 是 I 中的一个符号, 或者是表示 x 的一个子串的符号。例如, 如果 $x = abbabb, y = ababbbaa$, 那么 $[1:2][4:6]aa$ 是长度为 4 的 y 的一个表示, 符号 $[i:j]$ 表示 x 的子串 $a_i a_{i+1} \dots a_j$ 。[提示: 使用习题 9.33 的结论。]
- 9.35 试证明长度为 n 的字符串的压缩位置树至多有 $4n-2$ 个顶点。
- **9.36 设计一个 $O(n)$ 算法为长为 n 的字符串构造压缩位置树。
- 9.37 字符串 $x = a_1 a_2 \dots a_n$ 是 $y = b_1 b_2 \dots b_p$ 的一个子序列, 如果对某个 $i_1 < i_2 < \dots < i_n, a_1 a_2 \dots a_n = b_{i_1} b_{i_2} \dots b_{i_n}$ (即 x 是 y 中的 0 个或多个符号被删除后得到的字符串)。试设计一个 $O(|x| \cdot |y|)$ 的算法来找到 x 和 y 的最长公共子序列。
- 9.38 给定字符串 x 和 y , 设计一个算法以确定为实现 x 转换到 y 所需的最短的单个符号的插入和删除的操作序列。

研究型问题

- 9.39 习题 9.10 中的时间限制能够进一步改善吗?
- 9.40 为了确定正则表达式 α 表示的语言中的字符串是否 x 的子串, $O(|\alpha| \cdot |x|)$ 时间是否必要?
- 9.41 一个 k 磁头 2DPDA 能够在 $O(n^k)$ 时间内被一台 RAM 所模拟(习题 9.21)。那么每个上下文无关语言都能够被 2 磁头 2DPDA 接受吗? 如果能, 那么每个上下文无关语言都能够在 $O(n^2)$ 时间内被一台 RAM 所接受。
- 9.42 存在不能被 2DPDA 接受的上下文无关语言吗?
- 9.43 存在能够被 2NPDA 接受、但不能被 2DPDA 接受的语言吗?
- 9.44 习题 9.37 中的时间界限能够改善吗? 请查阅 Aho, Hirschberg 和 Ullman[1974] 基于决策树模型所得到的时间下界的结论。

文献与注释

Kleene[1956]研究了有穷自动机和正则表达式之间的等价性。Rabin 和 Scott[1959]研究了非

确定型有穷自动机,并证明了它们与确定型有穷自动机的等价性。正则表达式的模式匹配算法(算法 9.1)是 Thompson[1968]所提出的算法抽象。Ehrenfeucht 和 Zeiger[1974]讨论了 NDFA 作为模式规范设备的复杂度问题。一个字符串和另一个字符串的线性匹配算法(算法 9.3)来自于 Morris 和 Pratt[1970]。作为习题 9.21 的扩展, Cook[1971a]研究了 2DPDA 的线性模拟。Gray、Harrison 和 Ibarra[1967]研究了 2DPDA 的属性特征。Aho、Hopcroft 和 Ullman[1968]研究了双向非确定型 PDA 的 $O(n^3)$ 模拟算法(习题 9.24)。习题 9.15(b)来自 D. Chester 的研究,习题 9.19(c)来自于 V. Pratt。习题 9.19(c)的解决方案包含在 Knuth 和 Pratt[1971]的研究里。

9.5 节关于位置树和压缩位置树的材料来自于 Weiner[1973]。习题 9.20、习题 9.31 ~ 9.36 的解决方法及其几个相关应用的内容可在 Knuth[1973b]的研究中找到。Karp、Miller 和 Rosenberg[1972]的论文也包含了一些令人感兴趣的关于重复子串的模式匹配算法。关于 SWDC 匹配的习题 9.9 和 9.10 来自于 Fischer 和 Paterson[1974]的研究, Wagner 和 Fischer[1974]的研究则包含了习题 9.38 的解决方案。这篇论文和 Hirschberg[1973]阐述了习题 9.37 中的公共子序列问题的算法。

第 10 章 NP 完全问题

一个问题的计算量应该是多少才能将其评估为真正的难题？普遍的观点认为，如果一个问题不可能在少于指数时间内解决，那么这个问题应该认为是完全不易处理的。这个“分级模式”意味着有多项式时间界的算法问题是易解的。尽管诸如 2^n 这样的指数函数的增长比任何 n 的多项式的函数增长得快，但对于一些较小的 n 值， $O(2^n)$ 时间界的算法可能比多项式时间界的算法更有效。例如， 2^n 的值在 n 达到 59 之前都小于 n^{10} 。然而，一个指数函数的增长速度是爆炸性的，如果解决相应问题的所有算法都至少是指数时间复杂度的，那么就可以说它是不易解决的。

这一章将证明有一类问题，即非确定多项式时间完全问题（简称“NP 完全”问题），很可能只包含不易解决的问题。这类问题包括许多“典型的”组合数学问题，如旅行商问题、哈密顿回路问题、整数线性规划问题，以及所有已证明和这一类问题“等价的”问题。从这个意义上说，如果一个问题易于解决的，则所有的问题都是易于解决的。数学家和计算机科学家已经对这类问题中的许多研究了数十年，但对其中任何一个问题都没有找到多项式时间界的算法，很自然可推测不存在多项式时间的算法，因此，可认为这一类中的所有问题是不易解决的。

本章还将考虑第二类问题，叫做“多项式空间完全”问题，它们至少和 NP 完全问题一样难以解决，但还不能证明它们是不易解决的。第 11 章给出了一些目前我们能够证明是不易解决的问题。

10.1 非确定型图灵机问题

下面将看到，NP 完全问题理论背后的关键点是非确定型图灵机问题^①。我们已经讨论过非确定型有穷自动机，机器每一步移动都可能进入几个状态中的某一个。类似地，一个非确定型图灵机在每一步都有有穷数目的移动可供选择。如果对输入串 x ，至少存在一个移动序列导致机器到达一个接受的即时描述(ID)，则称输入串 x 可被机器接受。

给定输入 x ，我们能够把非确定型图灵机 M 想象成能并行执行所有可能的移动序列直到机器到达接受 ID 或者不可能进行任何移动为止。也就是说，在机器作了 i 次移动后，我们能够认为存在多个 M 的“副本”，每个副本表示一个 M 在 i 步移动之后的 ID。在第 $(i+1)$ 步移动时，如果在副本 C 的 IDC 下，图灵机对下一步移动有 j 种选择，则此时副本 C 将自身复制成 j 个副本。

这样，输入 x 时 M 可能的移动序列可以排成一棵 ID 树。树上从根到叶子的每条路径都表示一个可能的移动序列。如果 σ 是终止在接受 ID 的最短移动序列，那么，一旦完成 $|\sigma|$ 个移动， M 就停机且接受输入 x 。“用”在处理 x 上的时间是 σ 的长度。如果输入 x 没有移动序列导致 M 进入接受 ID，那么 M 拒绝 x ，且对处理 x 所用的时间不作定义。

认为 M 只“猜测”在序列 σ 中的移动，且验证 σ 确实终止在一个接受 ID 上，那么问题的处理常常会很方便。然而，一个确定型机器在正常情况下不可能事先猜测出一个到达接受 ID 的移动序列，因此一个针对 M 的确定型模拟必须按照某种顺序找出 M 在 x 上的所有可能的移动序列树，直到发现一个终止在接受 ID 的最短序列为止。如果没有移动序列导致 M 进入接受状态，那么除非有一个关于最短接受序列的长度边界有先验的限制，否则 M 的确定型模拟可能永远进行

① 不熟悉图灵机的读者可参阅 1.6 节。

下去。这样, 自然就会猜想: 在非确定型图灵机能够执行的任务中, 有一些是不能被确定型机器在同等时间或空间复杂度下完成的。然而, 一个重要的开放问题是: 是否有语言能被某固定时间或空间复杂度的非确定型图灵机接受, 但是不能被同样时间或空间复杂度的确定型图灵机所接受。

定义 一个 k 带非确定型图灵机 (简称 NDTM) M 是一个 7 元组 $(Q, T, I, \delta, b, q_0, q_r)$, 这里所有的元素都和普通的确定型图灵机的含义一样, 除了移动函数 δ 现是从 $Q \times T^k$ 到 $Q \times (T \times \{L, R, S\})^k$ 的子集的映射, 即, 给定一个状态和 k 带符号列表, δ 返回下一次移动可进行的有穷可选集合; 每个选择包含一个新的状态、 k 个新的磁带符号和 k 个磁头的 k 个移动。注意, NDTM M 可能选择这些移动中的任何一个, 但是它不可能从一个移动中选择状态而从另一个移动中选择新的磁带符号, 或者作任何其他移动组合。

NDTM 的即时描述 ID 的确切定义和确定型图灵机 (DTM) 一样。一个 NDTM $M = (Q, T, I, \delta, b, q_0, q_r)$ 确定当前的状态比如 q , 再查看 k 个磁头扫描的每个符号, 比如 X_1, X_2, \dots, X_k , 从而来作一次移动操作。然后从集合 $\delta(q, X_1, X_2, \dots, X_k)$ 选择某个元素 $(r, (Y_1, D_1), \dots, (Y_k, D_k))$, 该特殊元素规定了新状态为 r , 在第 i 个磁带上输出 Y_i 以代替 X_i , 并在 D_i 表明的方向上 ($1 \leq i \leq k$) 移动第 i 个磁头。如果 IDC 在选择了一个移动后变成 IDD, 记为 $C \vdash D$ (下标 M 经常被从下省略)。注意, NDTM M 中可能有几个 D , 满足 $C \vdash D$ 。但是, 如果 M 是确定性的, 则对每个 C 至多有一个这样的 D 。

对于某个 $k > 1$, 如果有 $C_1 \vdash C_2 \vdash \dots \vdash C_k$, 或者如果 $C_1 = C_k$, 则记为 $C_1 \vdash^* C_k$ 。如果 NDTM M 有 $(q_0 w, q_0, q_0, \dots, q_0) \vdash^* (\alpha_1, \alpha_2, \dots, \alpha_k)$, 这里 α_i (以及此后的 $\alpha_2, \alpha_3, \dots, \alpha_k$) 中的某个位置上有最后的状态符号 q_r , 则称 NDTM M 接受字符串 w 。 M 接受的语言记作 $L(M)$, 是 M 接受的所有字符串的集合。

例 10.1 设计一个 NDTM 接受形为 $10^{i_1} 10^{i_2} \dots 10^{i_k}$ 的字符串, 存在某个集合 $I \subseteq \{1, 2, \dots, k\}$, 有 $\sum_{j \in I} i_j = \sum_{j \notin I} i_j$ 。也就是说, 如果 w 代表的整数列表 i_1, i_2, \dots, i_k 能够分割成两个子列表, 其中一个子列表上的整数和等于另一个子列表上的整数和, 则 w 将被接受, 这个问题叫做分割问题。已经证明当整数用二进制编码, 问题的规模看成二进制整数列表的长度时^①, 该问题是 NP 完全的。

下面设计一个 3 带 NDTM M 来识别这个语言。它从左到右扫描它的输入磁带, 每次扫描一个 0 块到达 0^j , i_j 个 0 将不确定地添加到带 2 或带 3 上。当到达输入末尾时, NDTM 将检查它是否已经在带 2 上和带 3 上放置了相等数量的 0, 如果是, 则接受。这样, 如果按某种选择序列将各个 i_j 放到一个集合 (带 2) 或者另一个集合 (带 3) 中导致两集合内的整数和相等, 则 NDTM 将接受。而导致带 2 和带 3 上的整数和不相等的移动序列则不在考虑范围, 只要至少有一个选择序列导致接受即可。形式上记作

$$M = (\{q_0, q_1, \dots, q_5\}, \{0, 1, b, \$\}, \{0, 1\}, \delta, b, q_0, q_5),$$

其中移动函数 δ 如图 10-1 所示。

图 10-2 表明 NDTM 对输入 1010010 可能作的许多移动序列中的两个。第一个导致接受, 第二个则没有。既然至少有一个移动序列导致接受状态, 那么 NDTM 接受 1010010。 \square

① 因为元素可能有重复, 所以是一个列表而不是一个集合。

② 正如将看到的, 用来表示一个问题的编码是非常重要的。不难证明例 10.1 的 NDTM 接受语言的时间复杂度事实上为 $O_{TM}(n^2)$, 这里 n 是使用 DTM 时输入字符串的长度。然而, 如果用二进制对输入进行编码, 输入的长度为 i_1, i_2, \dots, i_k 诸对数之和, 同样的策略将得到一个 $O_{TM}(c^n)$ 算法, 这里的 n 是二进制输入的长度, $c > 1$ 。

定义 如果对于每个接受的长度为 n 的输入串有某个至多包括 $T(n)$ 次移动的序列导致 NDTM 进入接受状态, 则称 NDTM 的时间复杂度为 $T(n)$ 。如果对于每个接受的长度为 n 的输入字符串有某个移动序列导致 NDTM 进入接受状态, 在任何一个带上至多有 $S(n)$ 个不同的单元格被扫描, 则称 M 的空间复杂度为 $S(n)$ 。

例 10.2 例 10.1 的 NDTM 的时间复杂度为 $2n+2$ (最坏的情况是输入为 n 个 1), 空间复杂度为 $n+1$ 。其他合理的分割问题的编码也能得到相同的复杂度。例如, 用 $B(i)$ 表示整数 i 的二进制编码, 设

$$L_1 = \{ \#B(i_1)\#B(i_2)\cdots\#B(i_k) \mid \text{存在一个集合 } I \subseteq \{1, 2, \dots, k\}, \text{ 使得 } \sum_{j \in I} i_j = \sum_{j \notin I} i_j \}$$

这里 $\#$ 是一个特殊的标记符号。为了识别 L_1 , 可设计一个新的 NDTM₁, 它的操作类似于图 10-1 中的 NDTM M 。然而, M_1 不是将 0 块复制到带 2 或带 3, 而是将一个二进制数存储在带 2 和带 3。在输入带上遇到的每个新的二进制数被加到存储在带 2 或带 3 的数上。

状态	当前符号			(新的符号, 磁头移动)			新的状态	注 释
	带 1	带 2	带 3	带 1	带 2	带 3		
q_0	1	b	b	1, S	\$, R	\$, R	q_1	用 \$ 标记带 2 和 3 的左端, 然后转到状态 q_1
q_1	1	b	b	1, R	b , S	b , S	q_2	这里选择是否将下一块写到带 2 (q_2) 或带 3 (q_3)
				1, R	b , S	b , S	q_3	
q_2	0	b	b	0, R	0, R	b , S	q_2	复制 0 块到带 2。然后, 若在带 1 遇到 1, 则回到状态 q_1 ; 如果在带 1 遇到 b , 则转向状态 q_4 , 比较带 2 和 3 的长度
	1	b	b	1, S	b , S	b , S	q_1	
	b	b	b	b , S	b , L	b , L	q_4	
q_3	0	b	b	0, R	b , S	0, R	q_3	和状态 q_2 基本相同, 只是写到带 3 上
	1	b	b	1, S	b , S	b , S	q_1	
	b	b	b	b , S	b , L	b , L	q_4	
q_4	b	0	0	b , S	0, L	0, L	q_4	比较带 2 和带 3 的长度。
	b	\$	\$	b , S	\$, S	\$, S	q_5	
q_5								接受

图 10-1 NDTM 的移动, 每行代表一个选择

为了处理整数列表 i_1, i_2, \dots, i_k , M_1 将使用输入串 $x = \#B(i_1)\#B(i_2)\cdots\#B(i_k)$ 。为了处理同样的问题, M 将使用输入串 $w = 10^{i_1}10^{i_2}\cdots10^{i_k}$, 它的长度可能是 x 的指数级。这样, 尽管 M_1 在某种意义上比图 10-1 中的 NDTM 快指数系数级, 但由于 M_1 的输入已经相应缩短, 其时间和空间复杂度仍旧是 $O(n)$, 这里 n 是输入的长度。□

通过模拟可以证明任何 NDTM 接受的语言也能被 DTM 接受, 但是在时间复杂度上似乎要付出很多。还可以证明模拟的时间复杂度的最小上界是指数级的。也就是说, 如果 $T(n)$ 是一个合理的时间复杂度函数 (合理性是指“时间上可构造的”, 该概念将在第 11 章定义), 那么对于每个 $T(n)$ 时间界的 NDTM M , 可以找到一个常数 c 和一个 DTM M' 使得 $L(M) = L(M')$, 且 M' 的时间复杂度为 $O_{TM}(c^{T(n)})$ 。

上面结论的证明可以通过构造一个 DTMM' 并使用一个穷举算法来模拟 M 的实现。首先, 存在某个常量 d , 在任何情况下, M 的移动都不会多于 d 种选择。这样 M 的多达 $T(n)$ 个移动的序列能够用字母表 $\Sigma = \{0, 1, \dots, d-1\}$ 上长达 $T(n)$ 的字符串来表示。 M' 模拟 M 在长度为 n 的输

入 x 上的行为, 描述如下: M' 连续地按字典序产生长度至多 $T(n)$ 的 Σ^* 上的所有字符串, 这样的字符串不会多于 $(d+1)^{T(n)}$ 个。一旦产生了新串 w , M' 就开始模拟 σ_w , 这里 σ_w 是由 w 所代表的 M 的移动序列。如果 σ_w 引起 M 接受 x , 那么 M' 也接受 x 。如果 σ_w 不代表 M 的一个有效的移动序列, 或者如果 σ_w 没有引起 M 接受 x , 那么 M' 用 Σ^* 中的下一个字符串来重复这个过程。

M' 能够在时间 $O_{TM}(T(n))$ 内模拟 σ_w 。它至多使用 $O_{TM}(T(n))$ 时间产生每个字符串 w , 因此, NDTM 的整个模拟可能耗费时间 $O_{TM}(T(n)(d+1)^{T(n)})$, 至多是 $O_{TM}(c^{T(n)})$, 其中, c 是某个常量。模拟细节留作练习。

$(q_0 1010010, q_0, q_0)$	$(q_0 1010010, q_0, q_0)$
$\vdash(q_1 1010010, \$q_1, \$q_1)$	$\vdash(q_1 1010010, \$q_1, \$q_1)$
$\vdash(1q_2 010010, \$q_2, \$q_2)$	$\vdash(1q_3 010010, \$q_3, \$q_3)$
$\vdash(10q_2 10010, \$0q_2, \$q_2)$	$\vdash(10q_3 10010, \$q_3, \$0q_3)$
$\vdash(10q_1 10010, \$0q_1, \$q_1)$	$\vdash(10q_1 10010, \$q_1, \$0q_1)$
$\vdash(101q_3 0010, \$0q_3, \$q_3)$	$\vdash(101q_3 0010, \$q_3, \$0q_3)$
$\vdash(1010q_3 010, \$0q_3, \$0q_3)$	$\vdash(1010q_3 010, \$q_3, \$00q_3)$
$\vdash(10100q_3 10, \$0q_3, \$00q_3)$	$\vdash(10100q_3 10, \$q_3, \$000q_3)$
$\vdash(10100q_1 10, \$0q_1, \$00q_1)$	$\vdash(10100q_1 10, \$q_1, \$000q_1)$
$\vdash(101001q_2 0, \$0q_2, \$00q_2)$	$\vdash(101001q_3 0, \$q_3, \$000q_3)$
$\vdash(1010010q_2, \$00q_2, \$00q_2)$	$\vdash(1010010q_3, \$q_3, \$0000q_3)$
$\vdash(1010010q_4, \$0q_4 0, \$0q_4 0)$	$\vdash(1010010q_4, q_4 \$, \$0000q_4 0)$
$\vdash(1010010q_4, \$q_4 00, \$q_4 00)$	停机, 没有下一个 ID
$\vdash(1010010q_4, q_4 \$00, q_4 \$00)$	
$\vdash(1010010q_5, q_5 \$00, q_5 \$00)$	
接受	

图 10-2 对于一个 NDTM 的两个合法的移动序列

然而, 需要强调, 已知用 DTM 模拟 NDTM 的非平凡的时间下界。也就是说, 还不知道是否有语言能够被某个时间复杂度 $T(n)$ 的 NDTM 接受, 但是不能被同样时间复杂度的任何 DTM 接受。比较而言, 空间复杂度的相关问题更令人乐观。

定义 对于函数 $S(n)$, 如果存在一个 DTMM, 当给定长度为 n 的输入时, 它将在一个带的第 $S(n)$ 个格子上放置一个特殊的标记符, 而在任何带子占用的单元数不超过 $S(n)$ 个, 则说函数 $S(n)$ 是空间可构造的。大多数常用函数, 如多项式, 2^n , $n!$ 和 $\lceil n \log(n+1) \rceil$, 都是空间可构造的。

可以证明, 如果 $S(n)$ 是一个可构造的空间复杂度函数, M 是空间复杂度为 $S(n)$ 的 NDTM, 那么存在一个 DTMM' 使得 $L(M) = L(M')$, 且 M' 的空间复杂度为 $O(S^2(n))$ 。

M' 可以模拟 M 的策略是分治法的一个有趣应用。如果对于 k 带 NDTM $M = (Q, T, I, \delta, b, q_0, q_f)$, 其空间复杂度为可构造函数 $S(n)$, 那么存在常数 c , 使得当输入长为 n 的字符串时, M 需要进入的不同的 ID 数量至多为 $c^{S(n)}$ 。一个更加精确的界是 $\|Q\| \times (\|T\| + 1)^{kS(n)} \times (S(n))^k$, 其第一个因子表示状态的个数, 第二个因子限定了磁带格子内容的可能数量, 最后一个因子则限制了可能的磁头位置的数量。这样, 如果 $C_1 \stackrel{*}{\models} C_2$, 那么存在某个序列, 使得 M 在至多 $c^{S(n)}$ 次移动内从 IDC_1 转到 IDC_2 。可以通过对所有的 C_3 , 测试是否在至多 i 步移动内有 $C_1 \stackrel{*}{\models} C_3$ 和在至多

i 步移动内有 $C_3 \Vdash C_2$, 从而测试是否在至多 $2i$ 步移动内有 $C_1 \Vdash C_2$ 。

M' 之后的策略是通过如下的算法 10.1 来确定是否给定的初始 IDC_0 能够转到某个接受 ID。

算法 10.1 NDTM 的确定型模拟算法。

输入: 一个空间复杂度为 $S(n)$ 的 NDTM, 一个长度为 n 的输入串 w , 这里 $S(n)$ 是空间可构造的。

输出: 如果 $w \in L(M)$ 则输出“yes”; 否则输出“no”。

方法: 在图 10-3 里的递归过程 $TEST(C_1, C_2, i)$ 用于确定是否在至多 i 步内有 $C_1 \Vdash C_2$ 。如果是, 则返回值 true, 否则返回 false。在整个算法中, C_1 和 C_2 在任何带子上使用的单元数都不超过 $S(n)$ 。

完整的算法是对每个接受 IDC_i , 调用过程 $TEST(C_0, C_i, c^{S(n)})$, C_0 是输入为 w 的 M 的初始 ID。如果发现这些调用的任何一个值为 true, 则回答“yes”; 否则回答“no”。□

```

procedure TEST( $C_1, C_2, i$ ):
  if  $i = 1$  then
    if  $C_1 \vdash C_2$  or  $C_1 = C_2$  then return true
    else return false
  else
    begin
      for each ID  $C_3$  在任何的磁带上使用的格子数不多于  $S(n)$  个 do
        if TEST( $C_1, C_3, \lceil i/2 \rceil$ ) and TEST( $C_3, C_2, \lfloor i/2 \rfloor$ ) then
          return true;
    end
  end

```

图 10-3 过程 TEST

定理 10.1 如果 M 是一个复杂度为空间可构造函数 $S(n)$ 的 NDTM, 那么存在一个空间复杂度为 $O(S^2(n))$ 的 DTM M' , 使得 $L(M) = L(M')$ 。

证明: 证明是算法 10.1 的一个图灵机实现。由 2.5 节, 可以知道如何在 RAM 上模拟递归过程 TEST, 能够在图灵机上使用同样的堆栈策略。既然 TEST 的空间耗费参数是长度为 $O(S(n))$ 的 ID, 那么必须在一个磁带上安排对应于这个规模的栈帧; 而 $S(n)$ 是空间可构造的事实可确保我们进行这样的操作。检查 TEST 可以看到每次调用 TEST 时, 它的第三个参数基本上是 TEST 的调用过程的第三个参数的 $1/2$ 。因此, 可以证明任何时候栈帧的数量都不超过 $1 + \log \lceil c^{S(n)} \rceil$, 是 $O(S(n))$ 阶的。既然每帧有 $O(S(n))$ 单元被使用, 则图灵机中用于栈的单元总数是 $O(S^2(n))$ 。图灵机构造的余下细节留作练习。□

为了简化证明, 可将注意力集中在单带图灵机上。下面的引理指出如果愿意在计算时间上付出一些代价, 那么就可以将精力集中在单带图灵机上。

引理 10.1 如果 L 被一个时间复杂度 $T(n)$ 的 k 带 NDTM $M = (Q, T, I, \delta, b, q_0, q_r)$ 接受, 那么 L 可以被一个单带的时间复杂度为 $O(T^2(n))$ 的 NDTM 接受。

证明: 构造一个时间复杂度 $O(T^2(n))$ 的单带 NDTM M_1 以接受 L , 所用的策略是把 M_1 的磁带看做好像有 $2k$ “磁道”, 如图 10-4 所示, 也就是说, M_1 的磁带符号是 $2k$ 元组, 它的奇数号磁道的元素来自 T 的符号, 偶数号磁道的元素或者是空白 b 或者是一个特殊的标记符号 #。 k 个奇数号磁道相应于 M 的 k 个磁带, 每个偶数号磁道包含除了一个 # 符号外其余都是符号 b 。在 $2j$ 磁道上的符号 # 标记 M 在 j 磁带上的磁头位置, 它对应于 $2j-1$ 磁道。在图 10-4 中我们已经表明, 对每个 $j(1 \leq j \leq k)$, 第 j 个磁头扫描第 j 个磁带的 i_j 单元。

Tape 1 of M	X_{11}	X_{12}	...	X_{1i_1}	...
Head for tape 1	b	b	...	#	...
Tape 2 of M	X_{21}	X_{22}	...	X_{2i_2}	...
Head for tape 2	b	b	...	#	...
...					
Tape k of M	X_{k1}	X_{k2}	...	X_{ki_k}	...
Head for tape k	b	b	...	#	...

图 10-4 M_1 的磁带示意图

M_1 模拟 M 的一次移动如下。一开始, 假设 M_1 的磁头是在包含 M 的最左端磁头的单元上。

1. M_1 的磁头向右移动直到它经过偶数号磁道上的所有 k 个磁头标记。当它移动时, M_1 在自己的状态里记录 M 的每个磁头扫描的符号;

2. 在第 1 步中 M_1 的动作是确定性的, 现在 M_1 作一个非确定性的“分支”。具体地, 在 M 的状态(这个状态是 M_1 已经在自己的状态里记录下来的)和 M 扫描的磁带符号(这些符号是 M_1 刚确定的)的基础上, M_1 不确定地选择 M 的一个合法移动。对 M 在这种情况下的每个合法的移动, M_1 都有一个它可以进入的下一个状态;

3. 已经选择了 M 的一个移动进行模拟, M_1 根据选择的移动改变 M 的状态, 这个状态记录在自己的状态中。然后 M_1 向左移动磁头直到它再次经过 k 个磁头标记符。每当找到一个磁头标记符, 和选择的移动一致, M_1 在磁道上改变磁带符号, 且向左或向右移动磁头标记至多一格。

此时, M_1 已经模拟了 M 的一次移动。因为磁头至多在它的左端磁头标记右侧两格处, 所以标记可以被找到, 并使模拟过程循环重复。

既然 M_1 记录了 M 的状态, 那么, 如果 M 接受, 则 M_1 就可以接受。如果 M 接受长度为 n 的字符串 w , 则它用一个不多于 $T(n)$ 个移动的序列来完成接受工作。很明显, 一个包含 $T(n)$ 个移动的序列里, M 的磁头不可能到达多于 $T(n)$ 个不同的格, 所以 M_1 能够用自己的至多 $O(T(n))$ 个移动来模拟序列的一个移动。这样, M_1 通过包含至多 $O(T^2(n))$ 个移动的序列接受 w 。因此, M_1 接受语言 L , 且时间复杂度是 $O(T^2(n))$ 。□

推论 1 如果 k 带的时间复杂度 $T(n)$ 的 DTM 接受 L , 那么时间复杂度为 $O(T^2(n))$ 的单带 DTM 接受 L 。

证明: 在上面的证明中, 如果 M 是确定型的, 那么 M_1 也是确定型的。□

推论 2 如果空间复杂度为 $S(n)$ 的 k 带 NDTM 接受 L , 那么一定存在单带的空间复杂度 $S(n)$ 的 NDTM 接受 L 。

推论 3 如果存在空间复杂度为 $S(n)$ 的 k 带 DTM 接受 L , 那么一定存在单带的空间复杂度 $S(n)$ 的 DTM 接受 L 。

10.2 \mathcal{P} 类和 \mathcal{NP} 类

现在介绍两个重要的语言类:

定义 定义 $\mathcal{P}\text{-TIME}$ 是所有语言的集合, 这个集合中的每一语言都能够被多项式时间复杂度的 DTM 所接受, 也就是说:

$\mathcal{P}\text{-TIME} = \{L \mid \text{存在 DTM } M \text{ 和多项式 } p(n), \text{ 使得 } M \text{ 的时间复杂度为 } p(n), \text{ 且 } L(M) = L\}$ 。

定义 NP TIME 是所有语言的集合, 这个集合中的每一语言能够被多项式时间复杂度的 NDTM 接受。经常将 P TIME 和 NP TIME 分别简写成 P 和 NP 。

首先可以看到, 尽管已经用图灵机的术语给出 P 和 NP 的定义, 但是仍可以使用其他计算模型进行定义。直观地, 可以把 P 看做是在多项式时间内可识别的语言类。例如, 可以证明在对数代价模型下, 如果语言 L 关于图灵机的时间复杂度为 $T(n)$, 那么关于 RAM 或 RASP, 其时间复杂度处于 $k_1 T(n)$ 和 $k_2 T^*(n)$ 之间, 其中 k_1 和 k_2 为正常数。这样, 以对数代价为基准, L 能够被多项式时间复杂度的图灵机所接受当且仅当 L 在 RAM 或 RASP 上有多项式时间的算法。

还可以定义一个“非确定型”的 RAM 或 RASP, 方法是添加指令 CHOICE(L_1, L_2, \dots, L_k) 到 RAM 或 RASP 的指令列表里, 这个指令使得标记为 L_1, L_2, \dots, L_k 的语句中的一个语句非确定地选择和执行, 这样, 也能够在对数代价基准下, 用多项式时间界的非确定型 RAM 或者 RASP 定义 NP 类问题。

因此, 可以想像存在一个非确定型计算机, 诸如 RAM 或 RASP, 能够从一个给定的初始 ID 开始, 进行许多不同的可能的移动序列。这样的机器看来可以在多项式时间内识别许多语言, 而这些语言显然是不能在多项式时间内为确定型算法所识别的。当然, 在任何使用确定型机器 D 直接模拟非确定型机器 N 的尝试中, 由于 D 必须一直保持着记录 N 的大量副本, 所以为了执行所有可能的移动序列, D 所需要的时间比 N 的任何一个副本都多得多。能够从前一节的结果得出的最好结论是: 如果 L 是 NP 中的, 那么 L 能够被某个时间复杂度 $k^{p(n)}$ 的 DTM 所接受, k 为某个常数, p 为多项式, k 和 p 都依赖 L 。

另一方面, 目前还没有人能够证明一个语言在 NP 中而不在 P 中。也就是说, 还不知道 P 是否包含在 NP 中。然而, 可以证明某些语言是和 NP 中的语言一样“难”的, 在这种意义上, 如果有一个确定的多项式时间界的算法来识别这些语言中的一个, 就能够找到一个确定的多项式时间界的算法来识别 NP 中的任何语言。一般称这些语言是“NP 完全的”。

定义 如果 NP 中的语言 L_0 满足如下条件, 则定义 L_0 是非确定的多项式时间完全的 (简称 NP 完全的): 如果给定一个时间复杂度 $T(n) \geq n$ 的确定型算法来识别 L_0 , 那么对 NP 中的每个语言 L , 都能够有效地找到一个时间复杂度 $T(p_L(n))$ 的确定型算法来识别, 这里 p_L 是依赖于 L 的多项式。一般称 L 是可归约到 L_0 的。

证明语言 L_0 是 NP 完全的方法是证明 L_0 在 NP 中, 且每个 NP 中的语言 L 可以“多项式转换”到 L_0 。

定义 如果有一个多项式时间界的确定型图灵机 M , 它将字母表 L 里的每个串 w 转换成字母表 L_0 中的串 w_0 , 使 w 在 L 中当且仅当 w_0 在 L_0 中, 则定义语言 L 是多项式可转换到 L_0 的。

如果 L 可转换到 L_0 , 且 L_0 可被时间复杂度 $T(n) \geq n$ 的确定型算法 A 所接受, 那么能够确定输入 w 是否在 L 中, 先使用 M 将 w 转换到 w_0 , 然后使用算法 A 确定 w_0 是否在 L_0 中, 从而判断 w 是否在 L 中。如果 M 是 $p(n)$ 时间界的, 那么 $|w_0| \leq p(|w|)$ 。于是, 有一个算法确定 w 是否在 L 中, 所用时间为 $p(|w|) + T(p(|w|)) \leq T(2p(|w|))$ 。如果 T 是多项式的 (即 L_0 在 P 中), 那么接受 L 的算法将是多项式时间界的, 所以 L 也在 P 中。

事实上有一些作者给出如下的“ L_0 是 NP 完全”的定义: 如果语言 L_0 在 NP 中, 且每个 NP 中的语言都是多项式可转换到 L_0 的, 则称语言 L_0 是 NP 完全的。尽管不知道两个定义下的 NP 完全语言类是否不同, 但是这个定义看起来比前一个定义更严格。定义中“归约”的含义是指如果 M_0 是用于识别 NP 完全语言 L_0 的、确定的 $T(n)$ 时间界的图灵机, 那么 NP 中的每个语言都能够在 $T(p(n))$ 时间内被确定型图灵机识别 (其中 p 为某个多项式), 该图灵机可把 M_0 作为一个子程序调用 0 次或多次。定义“转换”是指 M_0 只用一次, 且只能以受限的方式调用。尽管我们应该采纳比较宽泛的定义, 但是我们所有的 NP 完全的证明却反映了比较狭义的定义。

关于这两个定义,以下一点是很明显的:如果有一个确定型的多项式时间界的识别 L_0 的算法,那么 NP 中的所有语言都能够在多项式时间内识别。这样,或者所有的 NP 完全语言都在 \mathcal{P} 中或者都不在。前者为真的当且仅当 $\mathcal{P} = NP$ 。遗憾的是,这里只能猜测 \mathcal{P} 是包含在 NP 中的。

10.3 语言和问题

上面把 \mathcal{P} 和 NP 定义为语言类,原因是双重的。第一,它简化了记法;第二,许多来自不同学科(如图论和数论)的问题经常可表示成语言识别问题。例如,考虑一个要求对每个实例回答“yes”或“no”的问题。可以把问题的每个实例编码为一个字符串,将初始问题形式化为一个语言识别问题,这样欲识别的语言就是包含所有回答为“yes”的问题实例的字符串表示。第9章讨论过这样的编码技术,那里使用语言识别的术语形式化表示许多模式匹配的问题。由于一个问题的时间复杂度可能依赖于所使用的编码,我们必须仔细挑选编码的方式。

为了使问题-语言之间的关系更明显,在此定义一些问题的公共“标准”表示。特别是可以进行以下假设。

1. 整数用十进制数表示;

2. n 个顶点的图用整数 $1, 2, \dots, n$ 表示图中的顶点,这些整数用十进制编码(如假设1),边用串 (i_1, i_2) 表示,这里 i_1 和 i_2 是表示边的顶点对的十进制表示;

3. 带 n 个命题变量的布尔表达式用字符串表示,在串中 $*$ 表示“and”, $+$ 表示“or”, \neg 表示“not”^①, 整数 $1, 2, \dots, n$ 表示命题变量。 $*$ 有时可以省略,如果需要的话,可以使用括号;这样就可以说一个问题是在 \mathcal{P} 或 NP 中当且仅当该问题的标准编码分别在 \mathcal{P} 或 NP 中。

例 10.3 考虑关于无向图的团问题。图 G 的 k 团是 G 中包含 k 个顶点的完全子图(在完全子图中,每一对不同顶点都相连一条边)。团问题是指:给定图 G 和整数 k , G 是否包含一个 k 团?

对于图 10-5 中的图 G , $k=3$ 时的团问题实例可以用如下字符串编码:

3(1, 2)(1, 4)(2, 3)(2, 4)(3, 4)(3, 5)(4, 5)

第一个整数表示 k 的值,然后是边所连接的顶点对,用 i 表示顶点 v_i , 即,按上面列出的次序,图中的边是 $(v_1, v_2), (v_1, v_4), \dots, (v_4, v_5)$ 。

表示团问题的语言 L 是如下形式的字符串集合。

$$k(i_1, j_1)(i_2, j_2) \dots (i_m, j_m)$$

满足带边 (i_r, j_r) ($1 \leq r \leq m$) 的图有一个 k 团。也可以使用其他语言来表示团问题。例如,常数 k 可能要求跟在图后面而不是放在图前面,或者用二进制整数代替十进制整数。然而,对于任何两个这样的语言,应该存在一个多项式 p , 使得一个语言中表示团问题实例的字符串 w 能够在时间 $p(|w|)$ 内转换成其他语言所表示的同样的问题实例的字符串。这样,就判断是否 \mathcal{P} 或 NP 中的成员而言,恰当地选择表示团问题的语言并不重要,只要使用“标准”形式的编码即可。□

团问题在 NP 中。一个非确定型图灵机能够先“猜测”哪 k 个顶点属于完全子图,然后在 $O(n^3)$ 步内验证这些顶点中的每对之间都有一条边,这里 n 是团问题的编码长度。因为 k 个顶点的所有子集可以被机器的独立副本并行地检查,所以非确定型图灵机的“强大”在此很明显地体现出来了。对于图 10-5 中的无向图,有三个 3 团,即 $\{v_1, v_2, v_4\}$, $\{v_2, v_3, v_4\}$ 和 $\{v_3, v_4, v_5\}$ 。

后面将看到团问题是 NP 完全的。目前为止,还没有在确定型多项式时间内解决团问题的已

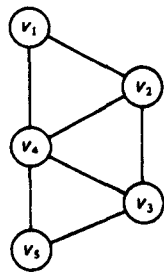


图 10-5 一个无向图

① 经常用 \neg 作为 $\neg(\alpha)$ 的简写。如果 α 只包括一个单个的文字(变量或者变量的补变量),那么括号可以省略。

知办法。

例 10.4 布尔表达式 $(p_1 + p_2) * p_3$ 可以用串 $(1+2)3$ 表示, 这里整数 i 表示变量 p_i 。考虑语言 L , L 包含所有表示可满足的布尔表达式的字符串(串中变量存在某个 0-1 赋值使表达式的值为 1)。后面将证明 L 是 NP 中的。一个接受 L 的非确定型算法通过“猜想”输入串中命题变量的一组满意的 0-1 赋值开始, 如果存在使布尔表达式值为 1 的变量赋值则猜想成功。然后, 用每个变量的值(0 或 1)去替换输入串中的对应变量。当用单个符号 0 或 1 替代命题变量的十进制表示时, 需要对串作移位, 然后计算结果表达式以验证其值是否为 1。

可以使用分解算法进行表达式的计算, 在和表达式长度成比例的时间内完成验证(见 Aho 和 Ullman[1972])。甚至, 不使用有效的分解算法, 读者也应该比较容易地在 $O(n^2)$ 步内计算出命题的值, 因此, 存在多项式时间复杂度的非确定型图灵机接受可满足的布尔表达式, 这样, 确定布尔表达式是否可满足的问题是在 NP 中的。因为“猜测”一个正确的解决方案实际上意味着并行测试所有可能的解决方案。这个问题将被证明是 NP 完全的。□

经常会对优化问题感兴趣, 如在图中寻找最大子团。许多这样的问题能够在多项式时间内转换为语言识别问题。例如, 图 G 的最大团可以用如下方法找到, 设 n 为 G 的表示长度, 对于从 1~ n 的每个 k , 确定是否有大小为 k 的团。一旦确定最大团的大小为 m , 就逐个地删除顶点直到移走顶点 v 会破坏所有余下的大小为 m 的团为止。然后考虑包含所有 G 中和 v 相邻的顶点构成的子图 G' , 递归地调用过程, 找到 G' 中大小为 $m-1$ 的团 C , G 的一个最大团就是 C 加上顶点 v 。

证明用上面的方法找到最大团的时间是 n 和 t 的多项式函数的问题, 留给读者作为练习, 其中, n 是 G 的表示长度, t 是判断 G 中是否有大小为 k 的团所用的时间。

一个形如“为某个命题 P 找到最大的 k 使得 $P(k)$ 为真”的优化问题经常能够用二分搜索法来解决(见 4.3 节), 这里 k 的可能取值关于问题描述长度 n 是指数数量级的。如果 $P(k)$ 为真, 意味着 $P(i)$ 为真($i < k$, k 在 0 和 c^n 之间取值, c 为常数), 那么使得 $P(k)$ 为真的最大的 k 能够用二分搜索法在经过 $\log c^n = n \log c$ 次测试命题 $P(k)$ 后找到。这里读者再次需要确信 k 的最优值能够在一定时间内找到, 这个时间是 n 和确定 $P(k)$ 值的最大时间的多项式函数。

这里把优化问题留给聪明的读者, 后面只考察 yes-no 判定问题。

10.4 可满足性问题的 NP 完全性

我们能够证明一个问题或者更准确地说, 一个问题的语言表示 L_0 是 NP 完全的。需要证明两点: L_0 在 NP 中, 以及 NP 中的每个语言是多项式可转换到 L_0 的。因为非确定型的“能力”, 证明给定的问题是 NP 中的通常比较容易。例 10.3 和 10.4 就是这一步的典型例子。最大的困难在于证明 NP 中的每个问题是多项式可转换到所给定的问题的。然而, 一旦证明了某个问题 L_0 是 NP 完全的, 就可以通过证明一个新的问题 L 在 NP 中, 且 L_0 是多项式可转换到 L 上, 从而证明 L 是 NP 完全的。

下面将证明布尔表达式的可满足性问题是 NP 完全的。然后再证明一些命题演算和图论中的基本问题是 NP 完全的, 证明是通过先证明它们在 NP 中, 然后再证明可满足性问题(或者一些别的已经证明是 NP 完全的问题)是多项式可转换到它们之上的, 从而证明这些问题是 NP 完全的。

阐述无向图上一些重要的 NP 完全问题时需要用到的定义如下。

定义 设 $G = (V, E)$ 表示一个无向图。

1. G 的顶点覆盖是 V 的子集 S , $S \subseteq V$, 满足 G 中每条边都与 S 中的某个顶点相关联。
2. G 的哈密顿回路是包含 V 中每个顶点的一条回路。
3. 如果对 G 的顶点存在一组称为“颜色”的整数 $1, 2, \dots, k$ 的赋值, 使得没有两个相邻的

顶点赋予相同颜色, 则称 G 是 k 可染色的。 G 的色彩编号是指使得 G 是 k 可染色的最小的整数 k 。

下面是后文中表示关于有向图的一些重要的 NP 完全问题时需要用到的定义。

定义 设 $G=(V, E)$ 是一个有向图。

1. 反馈顶点集合(a feedback vertex set)是一个顶点子集 $S, S \subseteq V$, 满足 G 的每个回路都包含 S 中的顶点。

2. 反馈边集合(a feedback edge set)是一个边的子集 $F, F \subseteq E$, 且 G 的每个回路都包含 F 中的边。

3. 有向哈密顿回路(a directed Hamilton circuit)是包含 V 的每个顶点的回路。

定理 10.2 如下的问题是 \mathcal{NP} 中的

1. (可满足性问题) 布尔表达式是可满足的吗?

2. (团问题) 无向图是否存在大小为 k 的团?

3. (顶点覆盖问题) 无向图是否存在大小为 k 的顶点覆盖?

4. (哈密顿回路问题) 无向图是否存在哈密顿回路?

5. (可染色性问题) 无向图是 k 可染色的吗?

6. (反馈顶点集问题) 有向图是否存在包含 k 个成员的反馈顶点集?

7. (反馈边集问题) 有向图是否存在包含 k 个成员的反馈边的集合?

8. (有向哈密顿回路问题) 有向图是否存在有向的哈密顿回路?

9. (集合覆盖问题) 给定一个集合族 S_1, S_2, \dots, S_n , 是否存在一个包含 k 个集合的子集合族 $S_{i_1}, S_{i_2}, \dots, S_{i_k}$, 使得

$$\bigcup_{j=1}^k S_{i_j} = \bigcup_{j=1}^n S_j?$$

10. (确切覆盖问题) 给定一个集合族 S_1, S_2, \dots, S_n , 是否存在包含两两不相交的子集合族的集合覆盖?

证明: 我们在例 10.4 和 10.3 中分别证明了问题 1 和 2 是 \mathcal{NP} 中的。类似地, 每一个其他问题都要求设计一个非确定型多项式时间界的图灵机(或者, 如果读者喜欢的话, 可以要求设计一个非确定型的 RAM 程序), 用这个图灵机去“猜测”一个解决方案并验证其是否真正地解决了问题。细节留作练习。 \square

现在证明 \mathcal{NP} 中的每个语言是多项式可转换到可满足问题上, 因而也就证明了布尔表达式的可满足问题是 NP 完全问题。

定理 10.3 确定布尔表达式是否可满足的问题是 NP 完全问题。

证明: 已经知道可满足性问题是 \mathcal{NP} 中的。这样只需证明 \mathcal{NP} 中的每个语言 L 是多项式可转换到可满足性问题上的即可。设 M 是一个接受 L 的, 多项式时间复杂度的非确定型图灵机, w 是 M 的一个输入。从 M 和 w 可以构造布尔表达式 w_0 , 使得 w_0 是可满足的当且仅当 M 接受 w 。证明的关键在于说明对每个 M 都存在多项式时间界的算法从布尔表达式 w 构造出 w_0 。这里多项式依赖于机器 M 。

根据引理 10.1, \mathcal{NP} 中的每个语言都可被一个多项式时间复杂度的非确定型单带图灵机所接受。这样, 可以假设 M 只有一个磁带, 且 M 有状态 q_1, q_2, \dots, q_t , 磁带符号是 X_1, X_2, \dots, X_m , 用 $p(n)$ 表示 M 的时间复杂度。

假设 M 的输入 w 的长度为 n 。如果 M 接受 w , 那么它在 $p(n)$ 个移动内完成。这样, 如果 M

⊙ 可用形如 $|i_1, i_2, \dots, i_m|$ 的串对有穷集合进行编码, i 为集合元素的十进制(或二进制)整数表示。如果在所有集合中有 n 个元素, 则用整数 j 表示第 j 个元素。

接受 w , 则至少存在一个 ID 序列 Q_0, Q_1, \dots, Q_q , 其中 Q_0 是初始 ID, $Q_{i-1} \vdash Q_i (1 \leq i \leq q)$, Q_q 是接受 ID, $q \leq p(n)$, 并且没有一个 ID 有多于 $p(n)$ 个磁带单元。

下面构造布尔表达式 w_0 来“模拟” M 所进入的一个 ID 序列。 w_0 中变量的真假赋值 (1 表示真, 0 表示假) 至多表示 M 的一个 ID 序列, 也可能不是合法的序列。布尔表达式 w_0 取值为 1 当且仅当变量赋值表示一个导致进入接受状态的 ID 序列 Q_0, Q_1, \dots, Q_q 。也就是说, w_0 可满足当且仅当 M 接受 w 。下面的变量是用在 w_0 中的命题变量, 其特定含义如下。

1. $C\langle i, j, t \rangle$ 是 1 当且仅当在时刻 t 时, M 的输入磁带的第 i 个单元包含磁带符号 X_j , 这里 $1 \leq i \leq p(n)$, $1 \leq j \leq m$, 且 $0 \leq t \leq p(n)$;
2. $S\langle k, t \rangle$ 是 1 当且仅当在时刻 t 时, M 处于状态 q_k , 这里 $1 \leq k \leq s$, 且 $0 \leq t \leq p(n)$;
3. $H\langle i, t \rangle$ 是 1 当且仅当在时刻 t 时, 磁头正在扫描磁带单元 i , 这里 $1 \leq i \leq p(n)$, 且 $0 \leq t \leq p(n)$ 。

这样, 有 $O(p^2(n))$ 个命题变量, 这些变量能够用至多 $c \log n$ 位二进制数表示 (c 为依赖 p 的常数)。在下面很容易维持这样一个假定: 每个命题变量能够用一个符号而不是 $c \log n$ 个符号来表示, 去掉因子 $c \log n$ 不会影响函数是否多项式时间界的问题。

用这些命题变量, 我们将根据 ID 序列 Q_0, Q_1, \dots, Q_q 构造一个布尔表达式 w_0 。在构造中, 利用谓词 $U(x_1, x_2, \dots, x_r)$, 该谓词在其参数 x_1, x_2, \dots, x_r 中恰好有一个为 1 时取值为 1。 U 能够表示成如下形式的布尔表达式:

$$U(x_1, x_2, \dots, x_r) = (x_1 + x_2 + \dots + x_r) \left(\prod_{\substack{i,j \\ i \neq j}} (\neg x_i + \neg x_j) \right) \quad (10-1)$$

式 (10-1) 的第一个因子保证至少一个 x_i 为 1, 余下的 $r(r-1)/2$ 个因子保证了没有两个 x_i 同时为 1。注意到当用形式化方法写出 U 时, $U(x_1, x_2, \dots, x_r)$ 的长度为 $O(r^2)^{\odot}$ 。

如果 M 接受 w , 那么有一个 M 在处理 w 时进入的 ID 接受序列: Q_0, Q_1, \dots, Q_q 。为了简化证明, 不失一般性, 假设 M 已经修改, 使得无论何时它到达接受状态时, 可以继续“运行”, 但既不会移动磁头也不会离开接受状态, 且序列中的每个 ID 都用空格延长至 $p(n)$ 。因此, 我们将 w_0 构造为下面七个表达式 A, B, \dots, G 的乘积, 以保证 Q_0, Q_1, \dots, Q_q 是一个接受的 ID 序列, 这里每个 Q_i 的长度为 $p(n)$ 且 $q = p(n)$ 。那么, 断言“ $Q_0, Q_1, \dots, Q_{p(n)}$ 是接受的 ID 序列”和如下断言等价:

1. 在每个 ID 中, 磁头恰好扫描一个单元;
2. 每个 ID 在每个磁带单元里恰好有一个磁带符号;
3. 每个 ID 恰好有一个状态;
4. 从一个 ID 到下一个 ID, 至多有一个磁带单元 (磁头扫描的单元) 被修改;
5. M 的移动函数允许连续的 ID 在状态、磁头位置和磁带单元内容上进行改变;
6. 第一个 ID 是初始 ID;
7. 最后一个 ID 中的状态是接受状态。

现在构造对应上面语句 1~7 的布尔表达式 $A \sim G$ 。

1. A 断言: 在每个时间单元, M 恰好扫描一个磁带单元。用 A_t 表示在时刻 t 时恰好有一个单元被扫描, 那么, $A = A_0 A_1 \dots A_{p(n)}$, 这里

$$A_t = U(H\langle 1, t \rangle, H\langle 2, t \rangle, \dots, H\langle p(n), t \rangle)$$

当缩写 U 被扩展时, A 的长度为 $O(p^3(n))$, 且能够在该时刻内写下来。

2. B 断言: 每个磁带单元在每个时刻单元恰好包含一个符号。用 B_i 断言第 i 个磁带单元在时

\odot 每个变量使用一个符号进行标记, 严格地说, 需要 $O(r^2 \log r)$, 至多需要 $O(r^3)$ 个符号。

刻 t 恰好包含一个符号。则

$$B = \prod_{i,t} B_{i,t}$$

这里

$$B_{i,t} = U(C\langle i, 1, t \rangle, C\langle i, 2, t \rangle, \dots, C\langle i, m, t \rangle)$$

因为磁带字母表的大小 m 只依赖于图灵机 M , 所以每个 $B_{i,t}$ 的长度独立于 n 。这样, B 的长度为 $O(p^2(n))$ 。

3. C 断言: 在每个时刻 t , M 在一个且仅在一个状态:

$$C = \prod_{0 \leq t \leq p(n)} U(S\langle 1, t \rangle, S\langle 2, t \rangle, \dots, S\langle s, t \rangle)$$

既然 M 的状态数 s 是一个常量, 那么 C 的长度为 $O(p(n))$ 。

4. D 断言: 在任何时刻 t 至多一个磁带单元的内容能够改变:

$$D = \prod_{i,j,t} [(C\langle i, j, t \rangle \equiv C\langle i, j, t+1 \rangle) + H\langle i, t \rangle]^\ominus$$

表达式 $(C\langle i, j, t \rangle \equiv C\langle i, j, t+1 \rangle) + H\langle i, t \rangle$ 断言:

a) 磁头在时刻 t 扫描单元 i , 或者

b) 在 $t+1$ 时刻, 第 j 个符号在单元 i 当且仅当在 t 时刻它在单元 i 。既然 A 和 B 声称磁头在时间 t 正扫描一个单元且单元 i 仅包含一个符号, 那么在时间 t 磁头或者正好扫描单元 i , 或者单元 i 的内容在时间 t 不作改变。当 \equiv 用相应的表达式替代时, D 的长度为 $O(p^2(n))$ 。

5. E 断言: M 的每个后继 ID 可以通过 M 的移动函数 δ 允许的一个转换从前一个 ID 获得。

用 E_{ijk} 表示如下断言之一:

a) 第 i 个单元在时刻 t 不包含符号 j ;

b) 磁头在时刻 t 不在扫描单元 i ;

c) M 在时刻 t 不在状态 k ; 或者

d) M 的下一个 ID 通过 M 的移动函数允许的转换从前一个 ID 获得。

则

$$E = \prod_{i,j,k,t} E_{ijk}$$

这里

$$E_{ijk} = \neg C\langle i, j, t \rangle + \neg H\langle i, t \rangle + \neg S\langle k, t \rangle + \sum_l [C\langle i, j_l, t+1 \rangle S\langle k_l, t+1 \rangle H\langle i_l, t+1 \rangle]$$

这里 l 的取值范围为: 当 M 在状态 q_i 扫描符号 X_j 时的所有可能的移动。也就是说, 对 $\delta(q_k, X_j)$ 中的每个 3 元组 (q, X, d) , 有一个 l 的对应值, 其中 $X_{j_l} = X$, $q_{k_l} = q$, 当 d 是 L, S 或 R 时 i_l 分别是 $i-1, i$ 或 $i+1$ 。这里 δ 是 M 的移动函数。因为 M 是非确定性的, 所以可能有多个这样的 3 元组 (q, X, d) , 但是数量是有限的。这样, E_{ijk} 是独立于 n 且长度有限的表达式, 因此, E 的长度为 $O(p^2(n))$ 。

6. F 断言: 初始条件被满足:

$$F = S\langle 1, 0 \rangle H\langle 1, 0 \rangle \prod_{1 \leq i \leq n} C\langle i, j_i, 0 \rangle \prod_{n < i \leq p(n)} C\langle i, 1, 0 \rangle$$

这里 $S\langle 1, 0 \rangle$ 表示在时间 $t=0$ 时, M 在状态 q_1 , q_1 是所选取的初始状态; $H\langle 1, 0 \rangle$ 表示在时间 $t=0$ 时, M 正在扫描左端的磁带单元; $\prod_{1 \leq i \leq n} C\langle i, j_i, 0 \rangle$ 断言磁带的前 n 个单元初始化为包含输入串 w ; $\prod_{n < i \leq p(n)} C\langle i, 1, 0 \rangle$ 断言余下的磁带单元初始化为空白单元。取 X_1 为空白符号, 显然, F 的长度为 $O(p(n))$ 。

7. G 断言: M 最终进入最后状态。因为已经要求 M 一旦到达最后状态就保持在最后状态, 因此有 $G = S\langle s, p(n) \rangle$ 。取 q_s 为最后状态。

⊖ 这里用 $x=y$ 表示 $xy + \bar{x}\bar{y}$, 即 x 当且仅当 y 。

布尔表达式 w_0 是乘积式 $ABCDEFG$ 。因为 w_0 的七个因子都需要最多 $O(p^3(n))$ 个符号, 那么 w_0 本身有 $O(p^3(n))$ 个符号。因为一直把命题变量看做单个符号, 而事实上它们是用长度为 $O(\log n)$ 的字符串表示, 因此实际上 $|w_0|$ 的界是 $O(p^3(n) \log n)$, 反过来它又受限于 $cn p^3(n)$ (其中 c 为某个常数)。这样, w_0 的长度是 w 长度的多项式函数。因此, 很明显, 给定 w 和多项式 p , w_0 能够在和与其长度成比例的时间内写下来。

给定一个接受的 ID 序列 Q_0, Q_1, \dots, Q_q , 能够容易地找到一组 0 和 1 赋值给命题变量 $C\langle i, j, t \rangle, S\langle k, t \rangle$ 和 $H\langle i, t \rangle$, 使得 w_0 的值为 1。相反, 给定使得 w_0 取值为 1 的一组变量赋值, 能够容易地找到一个接受的 ID 序列。这样, w_0 可满足当且仅当 M 接受 w 。

除了要求 M 接受的语言 L 在 \mathcal{NP} 中, 没有对 L 做任何限制。因此, 已经证明了任何 \mathcal{NP} 中的语言可多项式转换到布尔表达式的可满足性问题, 所以有结论: 可满足性问题是 NP 完全的。□

事实上定理 10.3 中的证明涉及到的内容比定理本身包含得多。如果一个布尔表达式是文字和的乘积, 那么该表达式是合取范式 (CNF)。这里对于某个变量 x , 文字是 x 或者 $\neg x$ 。例如, $(x_1 + x_2)(x_2 + \neg x_1 + x_3)$ 在 CNF 中, 而 $x_1 x_2 + x_3$ 则不在。在定理 10.3 中构造的表达式 w_0 实际上也在 CNF 中, 我们无需增加它的长度一个常数倍就能把它转换到 CNF。

推论 CNF 中的布尔表达式的可满足性问题是 NP 完全的。

证明: 完全可以证明, 在定理 10.3 的证明中定义的 A, \dots, G 中的每个表达式已经在 CNF 中, 或者能够通过使用布尔代数运算法则, 不用增加表达式的长度多于一个常量倍将其改为 CNF 范式。在公式 (10-1) 中定义的 U 已经在 CNF 中, 因此 A, B 和 C 在 CNF 中, 因为 F 和 G 是单个文字的乘积, 所以 F 和 G 显然在 CNF 中。

D 是形式为 $(x \equiv y) + z$ 的表达式乘积。如果替换 \equiv 符号, 则可得表达式

$$xy + \bar{x}\bar{y} + z \quad (10-2)$$

可以看出, 表达式 (10-2) 等价于

$$(x + \bar{y} + z)(\bar{x} + y + z) \quad (10-3)$$

用式 (10-3) 替代式 (10-2) 在 D 中的所有出现, 能够得到一个等价的 CNF 表达式, 且其长度至多两倍于原式的长度。

最后, 表达式 E 是表达式 E_{ijk} 的乘积, 既然每个 E_{ijk} 长度的界独立于 n , 那么每个 E_{ijk} 能够用长度独立于 n 的 CNF 表示。这样, E 转换到 CNF 将至多在长度上增加常数倍。

由此, 已经证明了表达式 w_0 能够转化为 CNF, 且得到的 CNF 的长度也只是比原来增加了常数倍。□

前面刚刚证明了 CNF 表达式的可满足性问题是 NP 完全的, 事实上, 即使在很严格的条件下, 布尔表达式的可满足性问题仍是 NP 完全的。如果一个表达式是至多 k 个文字的和的乘积, 则该表达式称为在 k 合取范式 (k -CNF) 中。 k 可满足问题确定是否 k -CNF 中的表达式可满足。对于 $k=1$ 或 2, 能够找到多项式时间的确定算法测试 k 可满足性问题, 而当 $k=3$ 时, 正如下面定理所描述的那样, 情况 (可以预见) 发生变化。

定理 10.4 3 可满足性问题是 NP 完全的。

证明: 我们将证明 CNF 表达式的可满足性可多项式转换到 3 可满足性。给定一个和的乘积, 用

$$(x_1 + x_2 + y_1)(x_3 + \bar{y}_1 + y_2)(x_4 + \bar{y}_2 + y_3) \cdots (x_{k-2} + \bar{y}_{k-4} + y_{k-3})(x_{k-1} + x_k + \bar{y}_{k-3}) \quad (10-4)$$

替代每个和式 $(x_1 + x_2 + \cdots + x_k)$, $k \geq 4$, y_1, y_2, \dots, y_{k-3} 为新变量。例如, 对于 $k=4$, 表达式 (10-4) 是 $(x_1 + x_2 + y_1)(x_3 + x_4 + \bar{y}_1)$ 。

某组新变量的赋值使得替代表达式的值为 1 当且仅当文字 x_1, x_2, \dots, x_k 中有一个值为 1。也就是说, 当且仅当初始的表达式值为 1。假设 $x_i = 1$, 那么当 $j \leq i-2$ 时, 设 y_j 为 1; 当 $j > i-2$,

设 y_i 为 0。此时替代表达式值为 1。反之, 假设对 y_i 的某组赋值使得结果表达式值为 1。如果 $y_1 = 0$, 那么 x_1 或者 x_2 必须为 1。如果 $y_{k-3} = 1$, 那么 x_{k-1} 或者 x_k 值必须为 1。如果 $y_1 = 1$ 且 $y_{k-3} = 0$, 那么对于某个 i , $1 \leq i \leq k-4$, $y_i = 1$ 且 $y_{i+1} = 0$, 这意味着 x_{i+2} 必须为 1。无论如何, 必须有一个 x_i 为 1。

每个替代表达式的长度是以被替代的表达式长度的常量倍为界的。事实上, 给定任意 CNF 表达式 E , 通过对每个和式运用上面的转换, 能够找到一个 3-CNF 表达式 E' , E' 是可满足的当且仅当初始表达式是可满足的。此外, 由于转换是直接应用的, 所以能够在和 E 的长度成比例的时间内找到 E' 。□

10.5 其他 NP 完全问题

通过直接或间接转换可满足性问题到定理 10.2 中提到的每个问题, 可证明它们也是 NP 完全的。图 10-6 的树结构说明了实际要进行的转换序列。在图 10-6 中, 如果 P 是 P' 的父结点, 那么证明 P 可多项式转换到 P' 。

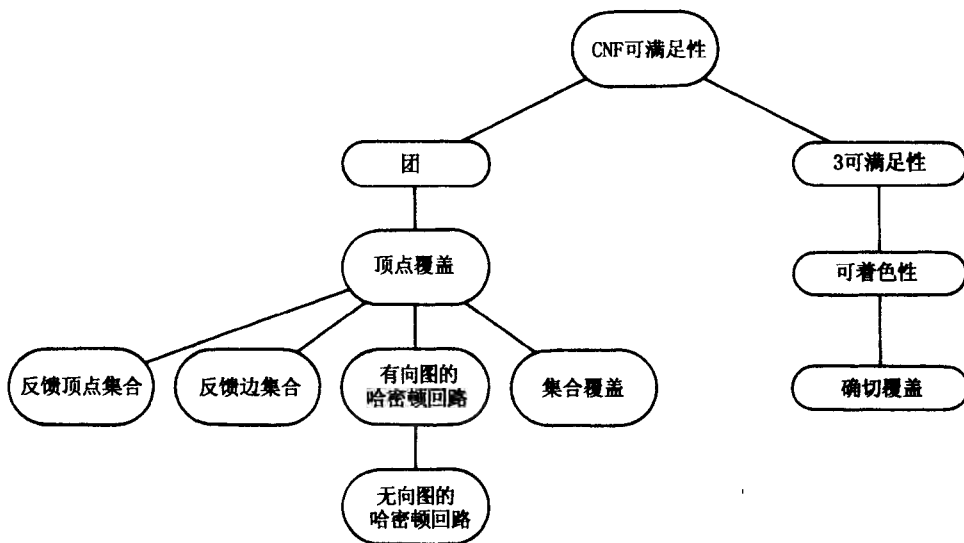


图 10-6 各种难题的转换顺序

定理 10.5 CNF 可满足性问题可多项式转换到团问题。因此团问题是 NP 完全的。

证明: 设 $F = F_1 F_2 \cdots F_q$ 是 CNF 中的表达式, 其中合取因子 F_i 的形式为 $(x_{i1} + x_{i2} + \cdots + x_{ik_i})$, 其中 x_{ij} 是文字。下面将构造无向图 $G = (V, E)$, 它的顶点是整数对 $[i, j]$ ($1 \leq i \leq q, 1 \leq j \leq k_i$)。整数对的第一个元素表示因子, 第二个元素表示因子中的文字。很自然, 图中的每个顶点对应特定因子的特定文字。

G 的边是顶点对 $([i, j], [k, l])$, 满足 $i \neq k, x_{ij} \neq \bar{x}_{kl}$ 。直观地, 如果 $[i, j]$ 和 $[k, l]$ 对应不同的因子, 且可以赋值给文字 x_{ij} 和 x_{kl} 中的变量, 使得两个文字的值均为 1 (即给 F_i 和 F_k 赋值为 1), 则在 G 中它们是相邻的。也就是说, $x_{ij} = x_{kl}$, 或者 x_{ij} 和 x_{kl} 是不同变量的求补形式或不求补形式。

在图 G 中的顶点数量明显小于 F 的长度, 边的数量则至多是顶点数的平方。这样, G 能够编码为一个串, 其长度限制在 F 长度的多项式边界内。更重要的是, 这样一个编码能够在 F 长度的多项式边界时间内计算。下面将证明 G 有大小为 q 的团当且仅当 F 是可满足的。因此, 给定一个团问题的多项式时间界的算法能够构造一个 CNF 可满足问题的多项式时间算法, 通过转换

表达式 F 到图 G , 使得 G 有大小为 q 的团当且仅当 F 是可满足的。接下来证明 G 有大小为 q 的团当且仅当 F 是可满足的。

当 假设 F 是可满足的, 那么存在对变量的 0-1 赋值, 使得 $F=1$ 。对于这个赋值, F 的每个因子值为 1, 每个因子 F_i 至少包含一个值为 1 的文字, 设 F_i 中的这个文字为 x_{im_i} 。

可以宣称顶点集合 $\{[i, m_i] \mid 1 \leq i \leq q\}$ 形成大小为 q 的团。否则存在 i 和 j ($i \neq j$), 使得顶点 $[i, m_i]$ 和 $[j, m_j]$ 之间没有边。根据 G 的边集定义, 这意味着 $x_{im_i} = \bar{x}_{jm_j}$ 。但是根据选择 x_{im_i} 的方法有 $x_{im_i} = x_{jm_j} = 1$, 所以这是不可能的。

仅当 假设 G 有大小为 q 的团, 因为第一个元素相同的两个顶点没有边连接, 所以团中的每个顶点的第一个元素肯定不同。因为团中恰好有 q 个顶点, 则团的顶点和 F 的因子之间一一对应。设团的顶点是 $[i, m_i], 1 \leq i \leq q, S_1 = \{y \mid x_{im_i} = y, \text{ 这里 } 1 \leq i \leq q, y \text{ 是一个变量}\}$, 且 $S_2 = \{y \mid x_{im_i} = \bar{y}, \text{ 这里 } 1 \leq i \leq q, y \text{ 是一个变量}\}$ 。也就是说, S_1 和 S_2 分别表示由团的顶点表示的不求补和求补变量的集合。那么, $S_1 \cap S_2 = \emptyset$, 否则在 $[s, m_s]$ 和 $[t, m_t]$ 之间将有一条边, 且 $x_{sm_s} = \bar{x}_{tm_t}$ 。通过设置 S_1 的变量为 1, S_2 的变量为 0, 可得每个 F_i 的值为 1。这样, F 是可满足的。 \square

例 10.5 考虑表达式 $F = (y_1 + \bar{y}_2)(y_2 + \bar{y}_3)(y_3 + \bar{y}_1)$, 其文字是

$$x_{11} = y_1, x_{21} = y_2, y_{31} = y_3$$

$$x_{12} = \bar{y}_2, x_{22} = \bar{y}_3, x_{32} = \bar{y}_1$$

根据定理 10.5 构造得到图 10-7。如图所示, 因为 $[1, 1]$ 和 $[1, 2]$ 的第一个元素相同, 所以没有边连接, 又因为 $x_{11} = y_1$ 且 $x_{32} = \bar{y}_1$, 所以 $[1, 1]$ 和 $[3, 2]$ 之间没有连接。 $[1, 1]$ 连接到图中其他三个顶点。

F 有三个因子, 在图 10-7 中碰巧有两个大小为 3 的团, 即 $\{[1, 1], [2, 1], [3, 1]\}$ 和 $\{[1, 2], [2, 2], [3, 2]\}$ 。在第一种情况下, 表示为三个文字 y_1, y_2, y_3 , 这些变量都没有求补, 且团相应于给 F 的赋值为 $y_1 = y_2 = y_3 = 1$ 。另一个团则相应于使 F 为 1 的另一个赋值, 也就是, $y_1 = y_2 = y_3 = 0$ 。 \square

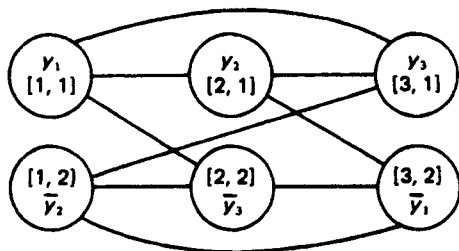


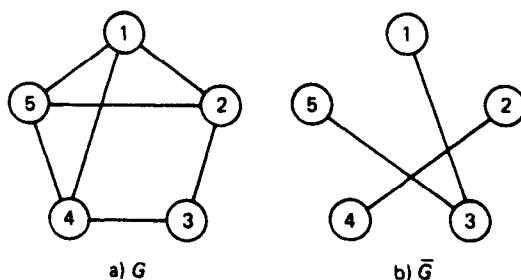
图 10-7 根据定理 10.5 构造的图

定理 10.6 团问题可多项式转换到顶点覆盖问题。因此顶点覆盖问题是 NP 完全的。

证明: 给定无向图 $G = (V, E)$, 考虑 G 的补图 $\bar{G} = (V, \bar{E})$, 这里 $\bar{E} = \{(v, w) \mid v, w \in V, v \neq w, \text{ 且 } (v, w) \notin E\}$ 。集合 $S \subseteq V$ 是 G 中的团当且仅当 $V - S$ 是 \bar{G} 的顶点覆盖。如果 S 是 G 的团, 在 \bar{G} 中没有边连接 S 中的两个顶点。这样, \bar{G} 中的每条边都至少附着在 $V - S$ 的一个顶点上, 意味着 $V - S$ 是 \bar{G} 的顶点覆盖。类似地, 如果 $V - S$ 是 \bar{G} 的顶点覆盖, 那么 \bar{G} 中的每条边至少附着在 $V - S$ 的一个顶点上。这样, \bar{G} 中没有边连接 S 中的两个顶点。因此 S 中的每对顶点均有 G 中的边连接, 则 S 是 G 的团。

为了弄清 G 是否有大小为 k 的团, 可构造 \bar{G} , 确定是否它有大小为 $\|V\| - k$ 的顶点覆盖。如果确实有, 则给定 $G = (V, E)$ 和 k 的标准表示, 能够找到 \bar{G} 和 $\|V\| - k$ 的一个表示, 且所需时间不超过 G 和 k 的表示长度的多项式时间。 \square

例 10.6 图 10-8a 的图 G 有大小为 3 的团 $\{1, 2, 5\}$ 和 $\{1, 4, 5\}$, 图 10-8b 中的 \bar{G} 有对应大小为 2 的顶点覆盖 $\{3, 4\}$ 和 $\{2, 3\}$; G 有大小为 2 的团 $\{2, 3\}$ 和 $\{3, 4\}$, \bar{G} 上有对应的大小为 3 的顶点覆盖 $\{1, 4, 5\}$ 和 $\{1, 2, 5\}$ 。 \square

图 10-8 a) 图 G b) G 的补图 \bar{G}

定理 10.7 顶点覆盖问题可多项式转换到反馈顶点集合问题。因此，反馈顶点集合问题是 NP 完全的。

证明：设 $G=(V, E)$ 是一个无向图，用两条有向边替换 G 中的每条边得到有向图 D 。具体地， $D=(V, E')$ ，其中 $E'=\{(v, w), (w, v) \mid (v, w) \in E\}$ 。因为 E 中的每条边用 D 中的一个回路代替，所以集合 $S \subseteq V$ 是 D 的反馈顶点集 (D 的每个回路都包含 S 中的一个顶点) 当且仅当 S 是 G 的顶点覆盖，而且很容易在多项式时间内从 G 中找到 D 的表示。 \square

定理 10.8 顶点覆盖问题可多项式转换到反馈边集合问题。因此反馈边集合问题是 NP 完全的。

证明：设 $G=(V, E)$ 是无向图， $D=(V \times \{0, 1\}, E')$ 是有向图，其中 E' 包括：

- 1) 对每个 $v \in V$ ，有边 $[v, 0] \rightarrow [v, 1]^{\ominus}$ ，且
- 2) 对于每条无向边 $(v, w) \in E$ ，有边 $[v, 1] \rightarrow [w, 0]$ 和 $[w, 1] \rightarrow [v, 0]$ 。

参阅图 10-9。

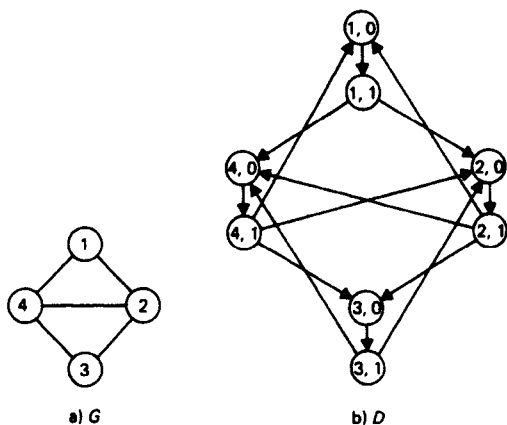


图 10-9 一个无向图 G 和对应于定理 10.8 中使用的有向图 D
 顶点覆盖 $\{2, 4\}$ 相应于反馈边集合 $\{(2, 0) \rightarrow (2, 1), (4, 0) \rightarrow (4, 1)\}$

设 $F \subseteq E'$ 是 D 的一个边集，对 D 中每个回路，该集合至少包含回路中的一条边，用边 $[w, 0] \rightarrow [w, 1]$ 代替 F 中的每条形式为 $[v, 1] \rightarrow [w, 0]$ 的边。记结果集为 F' ，则有 $\|F'\| \leq \|F\|$ ，且 F' 包含每个回路中至少一条边 (出自 $[w, 0]$ 的仅有边指向 $[w, 1]$ ，因此 $[w, 0] \rightarrow [w, 1]$ 是在包含 $[v, 1] \rightarrow [w, 0]$ 的任何回路中)。

\ominus 为了提高可读性，这里和本章的后续部分，将用 $x \rightarrow y$ 表示有向边 (x, y) 。

不失一般性, 假设对某个 k ,

$$F' = \{[v_i, 0] \rightarrow [v_i, 1] \mid 1 \leq i \leq k\}$$

那么对某个 $i (1 \leq i \leq k)$ 来说, D 的每个回路包含边 $[v_i, 0] \rightarrow [v_i, 1]$ 。然而, 如果 (x, y) 是 G 中的边, 那么 $[x, 1], [y, 0], [y, 1], [x, 0], [x, 1]$ 是 D 中的一个回路。这样 G 的每条边就附着在某个 v_i 上, $1 \leq i \leq k$, 因此 $\{v_1, v_2, \dots, v_k\}$ 是 G 的一个顶点覆盖。

相反, 我们可以方便地证明给定大小为 k 的顶点覆盖 S , 集合 $\{[v, 0] \rightarrow [v, 1] \mid v \in S\}$ 是 D 的反馈边集合。为了确定 G 是否有大小为 k 的顶点覆盖, 可以在多项式时间内构造 D , 且决定是否 D 包含有 k 个成员的反馈边集合。□

定理 10.9 顶点覆盖问题可多项式转换到有向图的哈密顿回路问题。因此有向图的哈密顿回路问题是 NP 完全的。

证明: 给定无向图 $G = (V, E)$ 和整数 k , 下面说明怎样在 $\|V\|$ 的多项式时间内构造一个有向图 $G_D = (V_D, E_D)$, 使得 G_D 有一个哈密顿回路当且仅当 V 的 k 顶点集合覆盖 G 的边。

设 $V = \{v_1, v_2, \dots, v_n\}$, 用 e_{ij} 表示边 $(v_i, v_j)^\ominus$, 引进新符号 a_1, a_2, \dots, a_k 。 G_D 的顶点包括相应于每个 a_i 的顶点加上 G 的每条边对应的四个顶点, 更精确地:

$$V_D = \{a_1, a_2, \dots, a_k\} \cup \{[v, e, b] \mid v \in V, e \in E, b \in \{0, 1\}, \text{且 } e \text{ 附着在 } v \text{ 上}\}$$

在形式化描述 G_D 的边之前, 先给出一个直观解释。如图 10-10 所示, G 的一条边 (v_i, v_j) 用一个带 4 个顶点的子图表示。

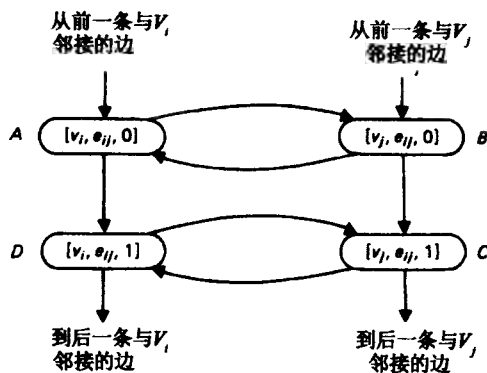


图 10-10 边 (v_i, v_j) 的表示

如果一条哈密顿回路从 A 进入表示边 e_{ij} 的子图, 那么它必须从 D 离开, 原因在于: 如果它从 C 离开, 则 $[v_j, e_{ij}, 0]$ 或者 $[v_i, e_{ij}, 1]$ 都不可能出现在回路中。从 A 到 D , 路径可以访问子图中的所有四个顶点, 也可能只访问左端的两个顶点。而在后一种情况下, 哈密顿回路必须在某点从 B 到 C 访问右端的两个顶点。

可以认为 G_D 包含 $\|V\|$ 个列表, 一个列表对应一个顶点。顶点 v 的列表包含所有按照特定顺序附着在 v 上的边。对 $1 \leq j \leq k$, 每个 a_j 有一条到顶点 $[v_i, e_{ij}, 0]$ 的边, 使得 e_{ij} 是 v_i 列表上的第一条边。如果 e_{ij} 是 v_i 列表的最后一条边, 则有一条从 $[v_i, e_{ij}, 1]$ 到每个 a_j 的边。如果在 v_i 的列表中边 e_{im} 紧跟在 e_{ij} 后, 则有一条边从 $[v_i, e_{ij}, 1]$ 到 $[v_i, e_{im}, 0]$ 。这些边加上图 10-10 包含的边形成集合 E_D 。现在证明 G_D 有一个哈密顿回路当且仅当 G 有覆盖所有边的 k 顶点集合。

当 假设顶点 v_1, v_2, \dots, v_k 覆盖 G 的所有边。记 $f_{i1}, f_{i2}, \dots, f_{i_{f_i}}$ 是附着在 v_i 上的边的列表

⊖ 注意 e_{ij} 和 e_{ji} 表示同一条边, 因此, 认为有着相同的符号。

($1 \leq i \leq k$)。考虑这样一个环路, 从顶点 a_1 到 $[v_1, f_{11}, 0]$, $[v_1, f_{11}, 1]$, $[v_1, f_{12}, 0]$, $[v_1, f_{12}, 1]$, \dots , $[v_1, f_{1i_1}, 0]$, $[v_1, f_{1i_1}, 1]$, 然后到 a_2 , 再到 $[v_2, f_{21}, 0]$, $[v_2, f_{21}, 1]$, \dots , 等等, 在通过 v_3, v_4, \dots, v_k 的所有边列表后, 最后回到 a_1 。除了不在与 v_1, v_2, \dots, v_k 的顶点相关的边列表中的顶点外, 这个环路遍历了 G_D 中的几乎所有顶点。既然 e 是附着在某个 v_i 上的, $1 \leq i \leq k$, 那么对于 G_D 的每对形如 $[v_j, e, 0]$, $[v_j, e, 1]$ ($j > k$) 的顶点, 可以扩展上述的环, 在环中用路径

$$[v_i, e, 0], [v_j, e, 0], [v_j, e, 1], [v_i, e, 1]$$

代替从 $[v_i, e, 0]$ 到 $[v_i, e, 1]$ 的边。对每个 v_j 和 e 如上修改的回路包含 G_D 的每个顶点, 因此, 是一个哈密顿回路。

仅当 假设 G_D 有一个哈密顿回路, 则环路能够拆成 k 段路径, 每条路径从顶点 a_i 开始, 在结点 a_j 结束。通过前面对表示每条边的四顶点子图中可能的路径的考虑 (如图 10-10), 可以看出存在顶点 v , 使得从 a_i 到 a_j 路径上的每个顶点的形式为 $[v, e, 0]$ 或 $[v, e, 1]$, 其中 e 是 G 中的边或者其形式为 $[w, e, 0]$ 或 $[w, e, 1]$, 其中 e 是边 (v, w) 。这样, 在从 a_i 到 a_j 的路径上的每个顶点的第一个元素或者是顶点 v , 或者是和 v 相邻的顶点。对于从 a_i 到 a_j 的路径, 可以使边和某个具体的顶点 v 相关联。对于 G_D 的每个顶点 $[w, e, b]$, e 一定是附着在和 G 中一条路径关联的 k 顶点中的一个顶点上。这样, k 个顶点形成 G 的一个顶点覆盖。因此, G_D 有一个哈密顿回路当且仅当 G 中存在 k 个顶点使得 G 的每条边附着在这 k 个顶点的一个顶点上。□

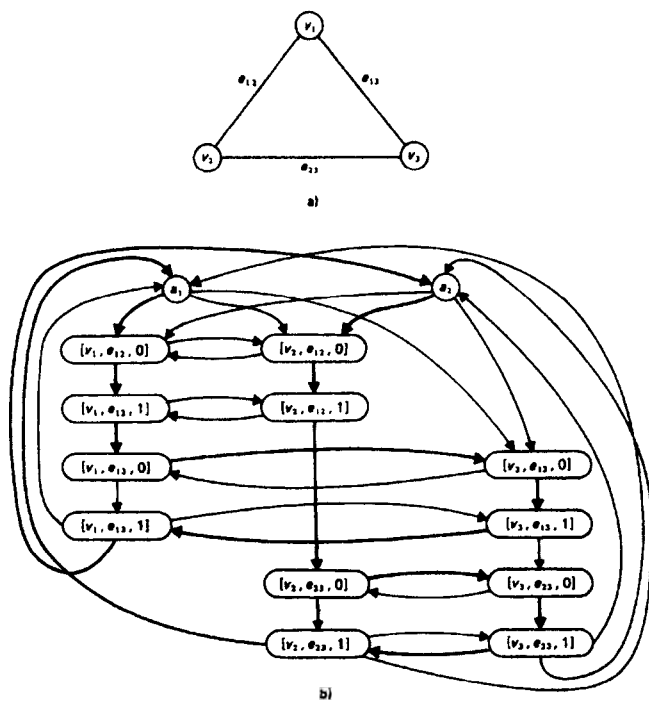


图 10-11 带哈密顿回路的图 a) 无向图 G b) 有向图 G_D

例 10.7 设 G 是有顶点 v_1, v_2 和 v_3 和边 e_{12}, e_{13} 和 e_{23} 的图, 如图 10-11a 所示。图 10-11b 是根据定理 10.9 构造的图 G_D , 图中粗连线表示基于顶点覆盖 $\{v_1, v_2\}$ 的 G_D 的一个哈密顿回路。□

定理 10.10 有向图的哈密顿回路问题可多项式转换到无向图的哈密顿回路问题。因此确定无向图中有否哈密顿回路的问题是 NP 完全的。

证明: 设 $G = (V, E)$ 是一个有向图, $U = (V \times \{0, 1, 2\}, E')$ 是一个无向图, 其中 E' 包含边

1. 对于 $v \in V$, E' 有边 $([v, 0], [v, 1])$;
2. 对于 $v \in V$, E' 有边 $([v, 1], [v, 2])$;
3. 当且仅当 $v \rightarrow w$ 是 E 中的一条边, E' 有边 $([v, 2], [w, 0])$ 。

注意到 V 中的每个顶点扩展成一条包含三个顶点的路径, 因为 U 中的哈密顿回路必须包括所有的顶点, 所以回路上顶点的最后一个元素必须按照 $0, 1, 2, 0, 1, 2, \dots$ 或者相反的顺序变化。可以假设其按前者变化, 其中类型 3 的边, 其第二个元素从 2 变到 0, 表示 G 中的有向哈密顿回路。反之, 通过使用 U 中的路径 $[v, 0], [v, 1], [v, 2], [w, 0]$ 替代每条边 $v \rightarrow w$, 则 G 中的一条有向哈密顿回路可以转化为 U 中的一条无向哈密顿回路。□

定理 10.11 顶点覆盖问题可多项式转换到集合覆盖问题。因此, 集合覆盖问题是 NP 完全的。

证明: 设 $G = (V, E)$ 是无向图, 对 $1 \leq i \leq \|V\|$, 设 S_i 是附着在 v_i 上的所有边的集合, 很显然, $S_{i_1}, S_{i_2}, \dots, S_{i_k}$ 是 S_i 的一个集合覆盖当且仅当 $\{v_{i_1}, v_{i_2}, \dots, v_{i_k}\}$ 是 G 的一个顶点覆盖。□

定理 10.12 3 可满足问题可多项式转换到可染色问题。

证明: 给定 3-CNF 中的带 n 个变量和 t 个因子的表达式 F , 证明如何在关于 $\text{Max}(n, t)$ 的多项式时间内构造带 $3n + t$ 个顶点的无向图 $G = (V, E)$, 使得 G 能够用 $n + 1$ 种颜色染色当且仅当 F 是可满足的。

设 x_1, x_2, \dots, x_n 和 F_1, F_2, \dots, F_t 分别是 F 的变量和因子, 引入新符号 v_1, v_2, \dots, v_n 。因为 3-CNF 中带 3 个或更少变量的任何表达式, 能够在表达式长度的线性时间内直接测试其可满足性, 而不用再转换到可染色性问题来检验, 因此不失一般性, 假设 $n \geq 4$ 。 G 中包含的顶点如下:

1. x_i, \bar{x}_i 和 $v_i, 1 \leq i \leq n$, 和
2. $F_i, 1 \leq i \leq t$ 。

G 中包含的边如下:

1. 对于 $i \neq j$, 所有边 (v_i, v_j) ;
2. 对于 $i \neq j$, 所有边 (v_i, x_j) 和 (v_i, \bar{x}_j) ;
3. 对 $1 \leq i \leq n$, 边 (x_i, \bar{x}_i) ;
4. 如果 x_i 不是因子 F_j 的项, 则为边 (x_i, F_j) ; 如果 \bar{x}_i 不是 F_j 的项, 则为边 (\bar{x}_i, F_j) 。

顶点 v_1, v_2, \dots, v_n 形成 n 个顶点的完全子图, 因此需要 n 种不同的颜色。每个 x_j 和 \bar{x}_j 均连接到每个 $v_i, i \neq j$, 因此, x_j 和 \bar{x}_j 除了可能和 v_j 颜色相同之外, 不可能和其他 v 顶点有同样的颜色。既然 x_j 和 \bar{x}_j 是相邻的, 它们也不可能是同一种颜色。所以, G 不可能用 $n + 1$ 种颜色染色, 除非 x_j 和 \bar{x}_j 的一个和 v_j 颜色相同, 而另一个是一种新的颜色, 我们称其为专色 (special color)。

考虑将 x_j 或 \bar{x}_j 中对应染成专色的那个顶点赋值为 0。现在考虑赋给 F_j 顶点的颜色。 F_j 至少和 $2n$ 个 x_i 和 \bar{x}_i 顶点中 $2n - 3$ 个顶点相邻。既然假设 $n \geq 4$, 那么对于每个 j 必存在一个 i , 使得 F_j 与 x_i 和 \bar{x}_i 都相邻。因为 x_i 或 \bar{x}_i 的一个是用专色染色的, F_j 就不可能用专色染色。

如果 F_j 包含某个文字 y , 而 \bar{y} 已经染为专色, 那么 F_j 就不能和任何与 y 同色的顶点相邻, 因此可以染成和 y 相同的颜色。否则, 需要一个新的颜色。这样, 所有的 F_j 顶点能够在无需增加颜色情况下染色当且仅当有一个对文字的专色赋值, 使得每个因子包含某个文字 y , 而 \bar{y} 染以专色。即当且仅当能够给变量赋值使得每个因子包含一个值为 1 的 y (\bar{y} 值为 0)。也可以说, 当且仅当 F 是可满足的。□

例 10.8 设 F 是 $(x_1 + x_2)(\bar{x}_1 + x_3)$ 。每个因子包含两项而不是三项, 这个变化并不改变定理 10.12 中图 G 的构造, 不过这使我们能够处理只涉及三个变量的表达式和能够用四种颜色染色的图(注意, 类似定理 10.12 对 2-CNF 也成立, 但这不是我们感兴趣的问题。因为对于 2-CNF 的可满足问题, 存在一个多项式时间界的算法进行验证, 所以 2-CNF 的可满足问题是多项式时间可转换到这里每个问题的)。结果如图 10-12 所示, 颜色标记成 A, B, C 和 S (专色), 其 4 色方案相应于变量赋值 $x_1 = x_2 = x_3 = 0$ 。□

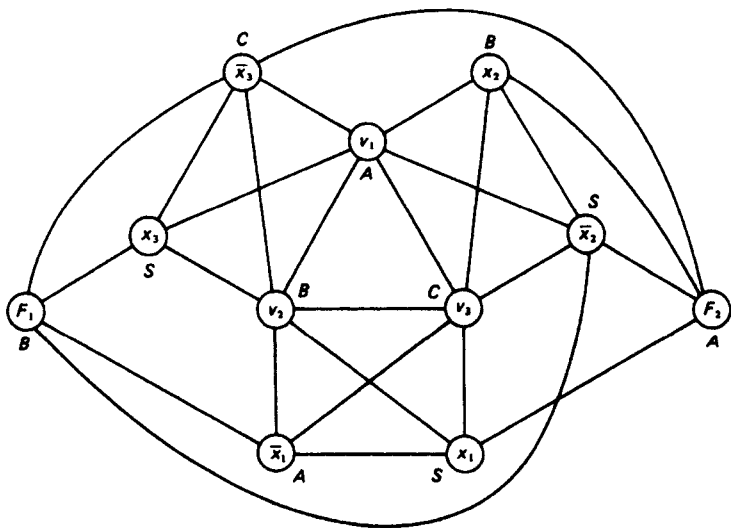


图 10-12 需要 4 色染色的图

定理 10.13 可染色问题可多项式转换到确切覆盖问题, 因此, 确切覆盖问题是 NP 完全的。

证明: 设 $G = (V, E)$ 是无向图, k 是整数, 下面构造集合, 集合的元素从 S 中选择

$$S = V \cup \{[e, i] \mid e \in E \text{ 且 } 1 \leq i \leq k\}$$

对每个 $v \in V$ 和 $1 \leq i \leq k$, 定义

$$S_v = \{v\} \cup \{[e, i] \mid e \text{ 附着在 } v \text{ 上}\} \quad (10-5)$$

对每条边 e 和 $1 \leq i \leq k$, 定义

$$T_{ei} = \{[e, i]\} \quad (10-6)$$

下面将图 G 的 k 染色问题关联到式(10-5)和式(10-6)定义的集合的确切覆盖问题。假设 G 有 k 染色方案, 其中顶点 v 染色成 c_v , 这里 c_v 可以取 $1 \sim k$ 的一个整数, 那么可以很容易地检查每个 v 对应的集合 S_v , 得到对于任何的 v , 满足 $[e, i] \notin S_v$ 的单元元素集合 T_{ei} 形成 S 的一个确切覆盖。在证明中注意到, 如果 $e = (v, w)$ 是一条边, 那么 $c_v \neq c_w$, 因此 $S_{c_v} \cap S_{c_w} = \emptyset$ 。当然, 如果 x 和 y 不相邻, 则 S_{c_x} 和 S_{c_y} 不相交, 且除非 T_{ei} 和所有其他被选集合都不相交, 否则, 不选择 T_{ei} 。

反之, 假设集合 S 有一个确切覆盖 \mathcal{S} 。因为每个 v 必须正好在 \mathcal{S} 的一个集合中, 那么对于每个 v , 就有唯一的 c_v 使得 S_{c_v} 在 \mathcal{S} 中。那么可以宣称每个顶点 v 可以染色 c_v 以获得 G 的 k 染色方案; 假设不是, 那么存在某条边 $e = (v, w)$, 使得 $c_v = c_w = c$ 。那么 S_{c_v} 和 S_{c_w} 都包含 $[e, c]$, 这样 \mathcal{S} 就不是一个确切覆盖, 与假设矛盾。□

10.6 多项式空间界问题

还有另一个普通的, 包含了 NP 难题的难题类, 称为 \mathcal{P} -SPACE, 是一个语言类, 能够被具

有多项式空间界的确定型图灵机所接受。

回忆一下,根据定理 10.1,如果 L 能被一个空间复杂度 $p(n)$ 的 NDTM 接受(p 为某个多项式),那么它就能被一个空间复杂度 $p^2(n)$ 的 DTM 接受,这样, \mathcal{P} -SPACE 也包括了所有被多项式空间界的 NDTM 接受的语言。因为每个多项式时间边界的 NDTM 是多项式空间界的,所以直接可得结论 \mathcal{NP} -TIME $\subseteq \mathcal{P}$ -SPACE(如图 10-13)。此外,因为 \mathcal{P} -TIME $\subseteq \mathcal{NP}$ -TIME $\subseteq \mathcal{P}$ -SPACE,如果 \mathcal{P} -TIME = \mathcal{P} -SPACE,则 \mathcal{P} -TIME = \mathcal{NP} -TIME,所以不可能出现 \mathcal{P} -SPACE 中的每个语言有确定型的多项式时间算法比 \mathcal{NP} -TIME 中的每个语言有确定型的多项式时间算法更不可能。

我们可以给出一个相当普通的语言 L ,它是 \mathcal{P} -SPACE 完全的,也就是说, L 在 \mathcal{P} -SPACE 中,如果给定一个 $T(n)$ 时间界的确定型 TM 接受 L ,那么对于 \mathcal{P} -SPACE 中的每个语言 L ,能够找到一个 $T(p_L(n))$ 时间界的确定型 TM 接受 L ,这里的 p_L 是依赖 L 的多项式。如果“确定型”换成“非确定型”,结论同样成立。这样,如果 L 在 \mathcal{P} -TIME 中,那么 \mathcal{P} -TIME = \mathcal{NP} -TIME = \mathcal{P} -SPACE;如果 L 在 \mathcal{NP} -TIME 中,那么 \mathcal{NP} -TIME = \mathcal{P} -SPACE。语言 L 是补集非空的正则表达式的集合。

引理 10.2 设 $M = (Q, T, I, \delta, b, q_0, q_f)$, 是 $p(n)$ 空间界的单带 DTM, p 是某个多项式,那么存在另一个多项式 p' , 如果 $x \in I^*$ 且 $|x| = n$, 就能够在时间 $p'(n)$ 内构造正则表达式 R_x , 使得 R_x 的补为空集当且仅当 M 不接受 x 。

证明: 证明和定理 10.3 非常相似,那里使用布尔函数作为描述图灵机计算的语言;这里则使用正则表达式语言。对正则表达式做一些精简,能够用正则表达式表示比正则表达式本身长得多的关于 ID 序列的事实。

猜想正则表达式 R_x 为表示 M 不接受 x 的 ID 序列。为此,使用字母表 $\Delta = T \cup \{[qX] \mid q \in Q \text{ 且 } X \in T\}$, $Q \times T$ 中的符号 $[qX]$ 表示当 M 在状态 q 时,输入磁带被磁头扫描的单元符号为 X 。如果 M 接受 x ,那么它通过一个移动序列完成接受操作,在这个序列中 M 至多使用 $p(|x|)$ 个磁带单元格。这样,可以用长度正好是 $p(n)$ 的符号串 w_i 表示 ID,符号串中除了一个符号之外,其他符号都来自 T ,剩下的符号形式是 $[qX]$ 。 M 的移动序列能够用串 $\#w_1\#w_2\#\dots\#w_k\#$ ($k \geq 1$) 表示,这里 $\#$ 是一个新的分隔符,每个 w_i 是 Δ^* 中的一个字符串,表示一个 ID, $|w_i| = p(n)$ 。必要时,为了使 w_i 长度为 $p(n)$,需要用空格填充 w_i 。

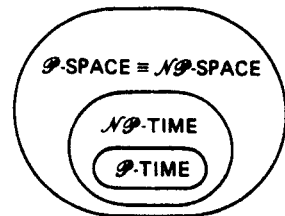


图 10-13 空间时间关系图

我们希望构造 R_x , 表示 $(\Delta \cup \{\#\})^*$ 中的那些字符串。它们表示 M 在输入为 x 时的不接受行为。这样, R_x 表示 $(\Delta \cup \{\#\})^*$ 中的所有字符串当且仅当 M 不接受 x 。只有当满足如下的四个条件中的一个或多个时, $(\Delta \cup \{\#\})^*$ 中的字符串 y 不表示 M 的一个接受计算。假设 $|x| = n$ 。

1. 对于某个 $k \geq 1$, y 的形式不是 $\#w_1\#w_2\#\dots\#w_k\#$, 其中, 对每个 i , $1 \leq i \leq k$, $|w_i| = p(n)$, 且 w_i 中除了一个符号, 其余符号都是 T 中的, 而其余符号形式是 $[qX]$;
2. 初始的 ID w_1 是错误的, 也就是说, y 不是以 $\#[q_0a_1]a_2\dots a_nb\#$ 开始, 这里 $x = a_1\dots a_n$, 且空白的数量是 $p(n) - n$;
3. M 从不进入最后状态, 也就是说, 没有形式为 $[q_fX]$ 的符号出现在 y 中;
4. y 有两个连续的 ID, 第二个 ID 不能通过 M 的一次移动从第一个 ID 得到。

我们用 $A + B + C + D$ 表示 R_x , 这里 A, \dots, D 分别表示以上条件 1, \dots , 4 所定义的正则表达式的集合。

给出一种正则表达式的简写方法是很有用的。如果有字母表 $S = \{c_1, c_2, \dots, c_m\}$, 用 S 表示正则表达式 $c_1 + c_2 + \dots + c_m$ 。另外, 使用 S^i 表示正则表达式 $(S)(S)\dots(S)$ (i 个), 也就是说, 用 S^i 表示字母表 S 上的长为 i 的字符串。注意到 S 表示的正则表达式的长度是 $2m - 1$, S^i 表示的

正则表达式的长度为 $i(2m+1)$ 。类似地, 如果 $S_1 - S_2 = \{d_1, d_2, \dots, d_i\}$, 那么用 $S_1 - S_2$ 表示正则表达式 $d_1 + d_2 + \dots + d_i$ 。

基于这些约定, 正则表达式 A 能够写成

$$\begin{aligned} A = & \Delta^* + \Delta^* \# \Delta^* + \Delta(\Delta + \#)^* + (\Delta + \#)^* \Delta \\ & + (\Delta + \#)^* \# T^* \# (\Delta + \#)^* \\ & + (\Delta + \#)^* \# \Delta^* (Q \times T) \Delta^* (Q \times T) \Delta^* \# (\Delta + \#)^* \\ & + (\Delta + \#)^* \# \Delta^{p(n)+1} \Delta^* \# (\Delta + \#)^* \\ & + A_0 + A_1 + \dots + A_{p(n)-1} \end{aligned} \quad (10-7)$$

A_i 的定义将在下文给出, 式(10-7)中的前 7 项分别表示:

- 1) 不带 # 的字符串;
- 2) 只带一个 # 的字符串;
- 3) 不是以 # 开头的字符串;
- 4) 不是以 # 结尾的字符串;
- 5) 在两个 # 之间没有 $Q \times T$ 符号的字符串;
- 6) 在两个 # 之间包含多于一个 $Q \times T$ 符号的字符串;
- 7) 在两个 # 之间包含多于 $p(n)$ 个 Δ 符号的字符串。

剩下的项 $A_i = (\Delta + \#)^* \# \Delta^i \# (\Delta + \#)^*$, 所有的 A_i 并在一起表示在两个 # 之间少于 $p(n)$ 个 Δ 符号的字符串的集合。

检查式(10-7)的正则表达式是否确实表示满足条件 1 的字符串, 留给读者完成。注意到前 6 项长度固定, 只依赖于 M ; 第 7 项的长度明显是 $O(p(n))$ 。项 A_i 长度为 $O(i)$, 所以式(10-7)所有项 A_i 长度的总和为 $O(p^2(n))$, 且比例常量只依赖于 M , 不依赖于 x 。此外, 容易验证表达式(10-7)能够容易地在 n 的多项式时间内写下来。

此时, 注意对于 B 、 C 和 D , 只需要写出满足条件 2、3 或 4 但不满足条件 1 的所有串的表达式即可, 也就是说, 只需要产生形式为 $\#w_1\# \dots \#w_k\#$ 的字符串, 这里 $|w_i| = p(n)$, 且 w_i 在 $T^*(Q \times T)T^*$ 中。然而, 为了简化问题, 这里将写出表达式以定义满足条件 2、3 和 4 但不满足条件 1 的串的集合, 也包括满足条件 2、3 或 4 以及条件 1 的一些字符串。因为根据表达式 A 的定义, 后一类中的串已经在 R_x 中, 所以它们在 B 、 C 和 D 中出现或不出现已经无关紧要了。

假设输入串为 $x = a_1 a_2 \dots a_n$, 将 B 表示为 $B_1 + B_2 + \dots + B_{p(n)}$, 其中对于 $1 < i \leq n$

$$\begin{aligned} B_1 &= \#(\Delta - [q_0 a_1]) \Delta^{p(n)-1} \# (\Delta + \#)^* \\ B_i &= \# \Delta^{i-1} (T - \{a_i\}) \Delta^{p(n)-i} \# (\Delta + \#)^* \end{aligned}$$

对于 $n < i \leq p(n)$,

$$B_i = \# \Delta^{i-1} (T - \{b\}) \Delta^{p(n)-i} \# (\Delta + \#)^*$$

这样, B_i 至少第 i 个输入符号不同于初始 ID。很明显, B 的长度是 $O(p^2(n))$, 比例常数只依赖于 M 。

把 C 写作 $((\Delta + \#) - (\{q_i\} \times T))^*$ 。显然, 这个表达式的长度只依赖于 M 。

最后, 按如下方式构造 D 。假定串 y 表示 M 的一个合法计算, 它的形式为 $\#w_1\#w_2\# \dots \#w_k\#$, 其中对每个 i , $|w_i| = p(n)$, w_{i+1} 是从 ID w_i 通过 M 的一次移动得到的 ID。观察给定 y 的任何三个连续符号 $c_1 c_2 c_3$, 能够唯一地确定符号, 称为 $f(c_1 c_2 c_3)$, 它一定在 c_2 的右边出现 $p(n) + 1$ 次。例如, 如果 c_1, c_2 和 c_3 都在 T 中, 那么 c_2 在下一个 ID 中可能不改变, 则 $f(c_1 c_2 c_3) = c_2$; 如果 c_1 和 c_2 在 T 中, $c_3 = Q \times T$ 中的 $[qX]$, 且 $\delta(q, X) = (p, X, L)$, 那么 $f(c_1 c_2 c_3) = [pc_2]$ 。定义 f 的其余规则, 包括: c_1 或 c_2 在 $Q \times T$ 中, 或者 c 中的一个为 # 符号时 f 的定义, 对读者来说这应该是很显然的, 留作练习。

因此, D 是对 $\Delta \cup \{\#\}$ 中所有 c_1 、 c_2 和 c_3 , 项 $Dc_1c_2c_3$ 的和, 其中

$$Dc_1c_2c_3 = (\Delta + \#)^* c_1c_2c_3 (\Delta + \#)^{p(n)-1} \bar{f}(c_1c_2c_3) (\Delta + \#)^*$$

其中 $\bar{f}(c_1c_2c_3) = \Delta \cup \{\#\} - f(c_1c_2c_3)$ 。我们看到 D 的长度是 $O(p(n))$, D 能够在 $O(p(n))$ 时间内写下来, 比例常量只依赖于 M 。

这样, 正则表达式 R_x 能够在 $p'(n)$ 时间内构造, 这里 p' 是 p^2 数量级的多项式。此外, R_x 表示 $(\Delta \cup \{\#\})^*$ 当且仅当 x 不被 M 所接受。□

引理 10.2 构造了一个正则表达式, 它的字母表 $\Delta \cup \{\#\}$ 的大小依赖于 M 。我们希望讨论“补集(关于它的字母表)是非空集合的正则表达式的语言”。既然一个语言一定有一个有穷字母表, 那么可在一个任意的字母表 Σ 上对一个正则表达式进行如下编码。

1. $+$, $*$, \emptyset , ϵ 和括号直接表示自身;
2. 字母表 Σ 的第 i 个符号(以任意次序)表示为 $\#w\#$, 这里 w 是 i 的十进制表示。

这样在任意的字母表上, 只要 17 个符号就可以编码任何的正则表达式, 设 $\Gamma = \{+, *, \emptyset, \epsilon, (,), \#, 0, 1, \dots, 9\}$ 是这个字母表。

定义语言 $L, \subseteq \Gamma^*$ 是正则表达式 R 的编码集合, 满足 R 的补集表示一个非空集合。如果 R 是字母表 Σ 上的, 那么补集也是关于 Σ 的。注意, 如果 R 的长度 $n \geq 2$, 那么它的编码长度至多为 $3n \log n$ 。

引理 10.3 L 在 \mathcal{P} -SPACE 中。

证明: 根据定理 10.1, 可以构造一个多项式空间界的 NDTM M 来接受 L 。设 R 是一个正则表达式, 它的编码作为输入字符串提供给 M 。根据定理 9.2, 给定长度为 n 的正则表达式, 构造一个非确定的、有至多 $2n$ 个状态的有穷自动机 A , A 接受正则表达式表示的语言。NDTM M 首先根据提供给 M 的 R 的编码建立自动机 A 。注意在定理 9.2 的构造中, A 的状态转移函数能够在不多于 $O(n^2)$ 的时间内构造。

假设 R 所表示集合的补集不空, 那么存在一个字符串 $x = a_1 \cdots a_m$, 而 x 不在 R 表示的语言里, M 可以一个符号一个符号地猜出 x 。 M 使用 $2n$ 位的数组保持跟踪记录状态集 S_i , S_i 是 A 在输入 $a_1, a_2, \dots, a_i (1 \leq i \leq m)$ 后可能处于的状态集。开始, S_0 是从 A 的初始状态仅通过 ϵ 转换可到达的状态集, 在 M 已经猜测了 a_1, a_2, \dots, a_i 之后, A 将处在某个状态集 S_i 中。当 M 猜测下一个输入符号 a_{i+1} 时, 它首先计算 $T_{i+1} = \{s' \mid s' \in \delta_A(s, a_{i+1}) \text{ 和 } s \in S_i\}$, 这里 δ_A 是 A 的状态转移函数; 然后对 T_{i+1} 中的每个 s' , 将 $\delta_A(s', \epsilon)$ 中的每个状态加到 T_{i+1} , 以计算 S_{i+1} (注意这和算法 9.1 类似)。

当 M 猜测 $a_1 a_2 \cdots a_m$ 时, 它将确定 S_m 没有包含 A 的最后状态。在找到一个没有最后状态的状态集时, M 接受它自己的输入串, 即正则表达式 R 的编码。这样, M 接受 R 的编码当且仅当 R 表示的集合的补非空。□

现在能够证明, 如果正则表达式所表示集合的补非空, 那么其编码组成的语言 L 是多项式空间完全的。

定理 10.14 如果 L 被时间复杂度 $T(n) \geq n$ 的 DTM 接受, 那么对于每个 $L \in \mathcal{P}$ -SPACE, 有一个多项式 p_L , 使得 L 在 $T(p_L(n))$ 时间内被接受。

证明: 设 M 是接受 L 的 DTM, 根据引理 10.2 和如下讨论, 可知有一个多项式时间界(如: 时间复杂度为 $p_1(n)$)的算法可以从输入串 x 构造正则表达式 R_x 。 R_x 定义在固定的 17 符号的字母表上, 长度一定不超过 $p_1(|x|)$ 。根据引理 10.2, R_x 的补集为空当且仅当 x 在 L 中。根据假设, 我们能够在 $T(|R_x|) = T(p_1(|x|))$ 步骤内测试 R_x 的补集是否空集。因此, 构造 R_x 的总时间是 $p_1(|x|)$, 测试 R_x 的时间是 $T(p_1(|x|))$ 。因为 $T(n) \geq n$, 所以选择 $p_L(n) = 2p_1(n)$ 足以证明定理。□

推论 L_i 在 \mathcal{P} -TIME 中, 当且仅当 \mathcal{P} -TIME = \mathcal{P} -SPACE。

证明: 如果 \mathcal{P} -TIME = \mathcal{P} -SPACE, 那么根据引理 10.3, L_i 在 \mathcal{P} -TIME 中。反之, 可以用 $T(n)$ 作为多项式从定理 10.14 得到。□

定理 10.15 如果 L_i 被时间复杂度 $T(n) \geq n$ 的 NDTM 接受, 那么对于每个 $L \in \mathcal{P}$ -SPACE, 有一个多项式 p_L 使得 L 在时间 $T(p_L(n))$ 内被一个 NDTM 所接受。

证明: 证明类似于定理 10.14。□

推论 L_i 在 \mathcal{NP} -TIME 中当且仅当 \mathcal{NP} -TIME = \mathcal{P} -SPACE。

习题

- 10.1 试给出图 10-1 中 NDTM 对输入 10101 的所有的合法移动序列, 判断 NDTM 是否接受这个输入。
- 10.2 试非形式化描述接受如下字符串集合的 NDTM 或非确定型简化 ALGOL 程序。
 - a) $10^i 10^j \dots 10^k$ 的字符串集合使得对 $1 \leq r < s \leq k$, 有 $i_r = i_s$;
 - b) xy 的字符串集合使得 x 和 y 在 $\{a, b\}^*$ 中, 且 x 是 y 的子串;
 - c) 和 (b) 一样, 但 x 是 y 的一个子序列。
- 10.3 试描述一个模拟非确定型 RAM 程序的 RAM 程序。
- 10.4 设 M 是 0 和 1 的 $m \times n$ 矩阵, 编写一个简化 ALGOL 程序, 求 M 的最小的 $s \times n$ 子矩阵 S , 如果 $M[i, j] = 1$, 那么对于某个 $1 \leq k \leq s$, 有 $S[k, j] = 1$, 并分析程序的时间复杂度。
- 10.5 非形式化地描述一个 DTM, 在其某条磁带的适当位置上放置一个标记符, 证明如下的函数是空间可构造的。
 - a) n^2
 - b) $n^3 - n^2 + 1$
 - c) 2^n
 - d) $n!$
- 10.6 试证明: 如果一个空间复杂度为 $S(n)$ 的单带 NDTM 有 s 个状态和 t 个磁带符号, 且接受长度为 n 的单词, 那么它在一个至多包含 $sS(n)t^{S(n)}$ 步移动的序列内接受这个单词。
- * 10.7 试说明怎样在 $O(n)$ 时间内, 对 RAM 计算一个布尔表达式。
- 10.8 试证明由不是重言式的布尔函数[⊖]组成的语言是 NP 完全的。
- 10.9 试证明子图同构问题(判断给定的无向图 G 是否和给定的无向图 G' 的某个子图同构)是 NP 完全的[提示: 使用团问题或者哈密顿回路问题是 NP 完全的事实]。
- 10.10 试证明背包问题(给定整数序列 $S = i_1, i_2, \dots, i_n$ 和整数 k , 是否存在 S 的子列, 其和恰好是 k ?)是 NP 完全的[提示: 可以使用确切覆盖问题是 NP 完全的事实]。
- * 10.11 试证明分割问题(习题 10.10 中 $\sum_{j=1}^n i_j = 2k$ 的背包问题)是 NP 完全的。
- 10.12 试证明旅行商问题(给定边为整数权值的图 G , 寻找一个回路, 它包括所有的顶点, 且边的权值之和至多为 k)是 NP 完全的。
- ** 10.13 试证明确定不含 * 的正则表达式(即只用运算符 + 和 \cdot)是否不能表示某个固定长度的所有字符串的问题是 NP 完全的[提示: 可以使用转换 CNF 可满足问题到正则表达式问题]。
- ** 10.14 试证明确定一个字母表 $\{0\}$ 上的正则表达式是否不表示 0^* 的问题是 NP 完全的。
- ** 10.15 设 $G = (V, E)$ 是一个无向图, v_1 和 v_2 是两个不同的顶点, v_1 和 v_2 的一个割集(cutset)是一个子集 $S \subseteq E$, 使得 v_1 和 v_2 之间的每条路径都包含 S 的一个元素。试证明确定是否一个图

⊖ 重言式是对于变量的所有取值, 其值均为 1 的布尔表达式。

有一个相应 v_1 和 v_2 的、大小为 k 的割集, 使得每部分都有相同顶点数的问题是 NP 完全的。

- ** 10.16 一维背包放置问题如下。 $G = (V, E)$ 是一个无向图, k 是一个正整数, 是否存在一个 V 的序列 v_1, v_2, \dots, v_n 使得

$$\sum_{(v_i, v_j) \in E} |i - j| \leq k?$$

试证明这个问题是 NP 完全的。

- ** 10.17 试证明, 即使 k 限制为 3, 顶点的最大度数限制为 4, 可染色问题依然是 NP 完全的。
 ** 10.18 试证明对于平面图习题 10.17 的结论。[提示: 先证明仅当 v_1 和 v_1' 同色、 v_2 和 v_2' 同色时, 图 10-14 的平面图才能用 3 色染色。将这个结论和习题 10.17 的证明结合起来。]
 * 10.19 试通过直接表示 NDTM 的计算而不是使用来自 3-CNF 可满足问题的转换来证明团问题是 NP 完全的。

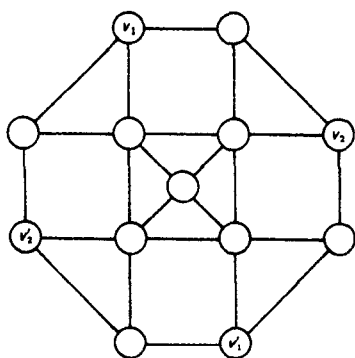


图 10-14

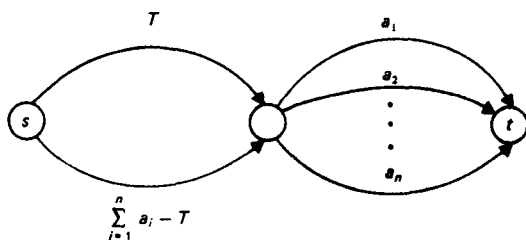


图 10-15

- * 10.20 考虑有两个指定顶点 s 和 t 的有向图, 每条边有一个整数“容量”, 可以通过反复寻找一条从 s 到 t 的路径, 并在边容量所允许的最大值范围内沿着路径增加流量, 从而构造最大流量图。试证明找到路径的最小集合, 在其上增加流量的问题是 NP 完全的[提示: 考虑如图 10-15 形式的图, 将问题和背包问题关联]。
- ** 10.21 相等执行时间的调度问题描述如下。给定作业集 $S = \{J_1, \dots, J_n\}$, 时间限制 t , 多个处理器 p , 和 S 上的偏序 $<$, 是否存在至多需要 t 个时间单元的 S 的调度, 即一个从 S 到 $\{1, 2, \dots, t\}$ 的映射使得至多 p 个作业映射到同一个整数上, 如果 $J < J'$, 则有 $f(J) < f(J')$? 证明这个问题是 NP 完全的。
- 10.22 试证明列在定理 10.2 中的问题是 \mathcal{AP} -TIME 的。
- 10.23 设 $p_1(x)$ 和 $p_2(x)$ 是多项式, 试证明存在一个多项式, 其对所有 x 的值都大于 $p_1(p_2(x))$ 。
- 10.24 如果 $\delta(q_k, x_j) = \{(q_3, X_4, L), (q_9, X_2, R)\}$, 请写出定理 10.3 证明中的 E_{ijk} 。
- ** 10.25 试给出一个测试 2 可满足问题的多项式算法。
- * 10.26 已经知道图同构问题的一些子问题如平面图的同构问题, 是比较容易处理的, 而其他子问题和一般的图同构问题处理难度是一样的。试证明对于带根的有向无回路图的同构问题和任意图的同构问题是一样复杂的。
- * 10.27 上下文敏感语言是一个可被空间复杂度为 $n+1$ 的 NDTM 接受的语言(叫做线性界自动机(linear-bounded automaton), 见 Hopcroft 和 Ullman [1969])。试证明确定给定的线性界自动机是否接受给定的输入问题是多项式空间完全的, 也就是说, 判断一个字符串是否属

于上下文敏感语言的问题是 \mathcal{P} -SPACE 完全的。

- 10.28 试证明确定两个正则表达式是否等价的问题是多项式空间完全的。
- **10.29 描述一个算法接受 L , L 是正则表达式 R 的编码集合, 其中 R 是能够被一个空间为线性界的 DTM 所实现的非空集合。

研究性问题

- 10.30 显然, 一个公开的问题是否 \mathcal{P} -TIME = \mathcal{NP} -TIME 或 \mathcal{NP} -TIME = \mathcal{P} -SPACE。考虑到已有的寻找 NP 完全问题的多项式时间算法的大量工作, 这个问题至少是和一些经典的数学问题如费马猜想(在对于 $n \geq 3$ 的整数范围内, $x^n + y^n = z^n$ 是否可解?) 或四色问题等一样困难。
- 10.31 如果不能解决习题 10.30 中的问题, 则还有一点更令人感兴趣是得到一个非平凡的函数 $T(n)$, 使得 \mathcal{NP} -TIME 中的每个语言的(确定型)时间复杂度为 $T(n)$ 。目前, 甚至 $T(n) = 2^n$ 的情况还没有被证明。

文献与注释

Hopcroft 和 Ullman[1969]给出了关于非确定的图灵机更多的信息。定理 10.1 将确定型和非确定型 TM 的空间复杂度相关联, 这是 Savitch[1970]的研究结果。

可满足性问题是 NP 完全的关键理论是 Cook 提出的[1971b]。Karp[1972]列举了大量典型的 NP 完全问题, 他给出了这一概念的重要性的清晰说明。从那时起, 许多科学家不断添加、扩充已有的 NP 完全问题族, 如: Sahni[1972], Sethi[1973], Ullman[1973], Rounds[1973], Ibarra 和 Sahni[1973], Hunt 和 Rosenkrantz[1974], Garey、Johnson 和 Stockmeyer[1974], Bruno 和 Sethi[1974]等等。令人感兴趣的是在 Cook 的开创性论文之前, 有几位学者证明了一些 NP 完全问题是多项式相关的, 但没有认识到此类问题的范围。例如, Dantzig、Blattner 和 Rao[1966]等人用负的权值把旅行商问题与最短路径问题相关联; Divetti 和 Grasselli[1968]证明了集合覆盖问题和反馈边集合问题的关系。

Meyer 和 Stockmeyer[1972]首先考虑了 \mathcal{P} -SPACE 完全问题。Savitch[1971]的研究结果提示了第一种空间完全语言的存在, 该论文定义了一种 $\log n$ 空间完全的语言(可解迷宫的集合)。Jones[1973]和 Stockmeyer 与 Meyer[1973]论述了某种限制形式的问题间的多项式时间可归约性。Book[1972 和 1974]证明了某些复杂问题类是不相容的。

习题 10.8 和 10.9 来自于 Cook[1971b]的研究结果, 习题 10.10 ~ 10.12 来自于 Karp[1972]的研究结果, 习题 10.13 来自 Stockmeyer 和 Meyer[1973]以及 Hunt[1973a], 习题 10.14 和 10.28 来自于 Stockmeyer 和 Meyer[1973], 习题 10.15 ~ 10.18 来自于 Garey、Johnson 和 Stockmeyer[1974], 图 10-14 根据 M. Fischer 的建议做了一些改进。平面 3 可染色问题是 NP 完全的证明出自 Stockmeyer[1973], 习题 10.20 来自 Even[1973], 习题 10.21 来自 Ullman[1973], 习题 10.25 来自 Cook[1971b], 习题 10.27 需要完成的推导来自 Karp[1972]。S. Even 给出了关于如何证明定理 10.13 的建议。

第 11 章 一些可证难的问题

本章将证明两类扩展正则表达式的补的空性问题是难解的, 即任何解决上述问题的算法至少需要指数阶的时间。对其中的一类, 本章给出一个实质上比指数阶大的下界。我们将特别说明该问题需要的时间多于

$$2^{2^{\dots^{2^n}}}$$

其中 2 的个数是有穷的。在证明下界之前, 先考虑表示“允许计算进行的时间和占用的空间越多, 则其能识别的语言就越多”的层次性结果。

11.1 复杂度层次

第 10 章说明了一些问题是非确定多项式时间完全或多项式空间完全的。为了证明一个特定的问题是完全的, 我们说明了需要使用该特定问题表示任何一个 NP -TIME 和 P -SPACE 中的问题。证明的技术基本是一种模拟。例如, 我们说明了布尔表达式的可满足性问题是 NP 完全的, 揭示了正则表达式补的空性问题是多项式空间完全的, 通过直接在这些问题的计算实例中嵌入图灵机计算来得到这些结果。我们可将其他问题归约到某类问题中已知是完全的问题, 来说明前者是完全的。因此, 也表明了 NP -TIME 和 P -SPACE 中都有“最难的问题”, 其复杂度至少和该类中的任何问题一样大。

然而, 引人注目的现实是, 目前还没有人发现一个问题在 NP -TIME 或 P -SPACE 中, 但不在 P -TIME 中。更进一步, 第 10 章的技术似乎仅可以说明一个问题至少和某一类中的其他问题是一样难的。为了实际证明一个问题不在 P -TIME 中, 需要一种技术说明至少存在一种语言不能在多项式时间内被任何确定型图灵机所接受。可使用的主要技术是对角线方法。尽管该技术不足以证明 P -TIME \neq NP -TIME, 但可用来建立确定型和非确定型图灵机的空间和时间复杂度的层次性。每个层次性定理都有如下形式: 给定两个“良好表现”的函数 $f(n)$ 和 $g(n)$, 其中 $f(n)$ 比 $g(n)$ 增长得“快”, 则存在一个复杂度为 $f(n)$ 但不为 $g(n)$ 的语言。

11.2 确定型图灵机的空间层次

本节证明一个关于确定性图灵机的空间层次性定理。对于时间和非确定型图灵机也有类似的层次性, 但和这里的目的需求不符, 因此留做练习。回忆一下引理 10.1 的推论 3, 如果语言 L 被一台空间复杂度为 $S(n)$ 的 k 带 DTM 所接受, 则 L 也能被空间复杂度为 $S(n)$ 的单带 DTM 所接受。因此讨论可限制在单带 DTM 上。

为了得到空间层次性结果, 需要对 DTM 进行仿真, 即对 DTM 指派一个顺序, 使得对每个非负整数 i 存在唯一的与 i 相随的 DTM。进一步要求每个 DTM 在仿真中出现得无限频繁。这里仅考虑仿真输入字母表为 $\{0, 1\}$ 的单带 DTM, 因为这就是所需要的一切。对所有图灵机的仿真方法是一种简单扩展。

不失一般性, 对单带 DTM 的表示可作如下假定。

1. 对于某个 s , 状态名为 q_1, q_2, \dots, q_s , 其中 q_1 为初始状态, q_s 为接受状态;
2. 输入字母表为 $\{0, 1\}$;

3. 对于某个 t , 带字母表为 $\{X_1, X_2, \dots, X_t\}$, 其中 $X_1 = b, X_2 = 0, X_3 = 1$;

4. 下一移动函数 δ 是形为 $(q_i, X_j, q_k, X_l, D_m)$ 的 5 元组列表, 含义为 $\delta(q_i, X_j) = (q_k, X_l, D_m)$, 其中 q_i 和 q_k 是状态, X_j 和 X_l 是带符号, D_m 是方向, 依 m 的值 0, 1 或 2 分别表示 L, R 或 S 。假定这个 5 元组列表由字符串 $10^i 10^j 10^k 10^l 10^m 1$ 编码。

5. 图灵机本身的编码是以任意次序将每个 5 元组的下一移动函数的编码连接在一起构成的。如果需要, 可在字符串的前面放置作为前缀的若干个 1。最后的结果是以 1 开始, 由 0 和 1 组成的字符串, 可解释为一个整数。

任何不能被解码的整数必定是表示了一个移动函数为空的平凡图灵机。在仿真中, 每台单带图灵机是无限频繁出现的, 因为给定一台 DTM, 可以随意增加前缀 1 而得到表示同一 5 元组集合的越来越大的整数。

现在可以构造一台 4 带图灵机 M_0 , 其将输入串 x 看成一台单带 DTMM 的编码, 同时也作为 M 的输入。我们将设计图灵机 M_0 , 使得对每台给定复杂度的 DTMM 至少存在一个输入串, 该输入串能被 M_0 接受但不能被 M 接受, 或者相反。 M_0 具备的能力之一是, 给定了一台图灵机的规范, 就可以仿真它。对某个函数 S_1 , 图灵机 M_0 可确定图灵机 M 是否使用不多于 $S_1(|x|)$ 个单元来接受输入串 x 。如果 M 在空间 $S_1(|x|)$ 内接受 x , 则 M_0 不接受 x ; 否则 M_0 接受 x 。因此, 对所有 i , 或者 M_0 不同意第 i 台 DTM 对于输入 x 的行为, 这里 x 是 i 的二进制表示, 或者第 i 台 DTM 对于输入 x 使用的单元数多于 $S_1(|x|)$ 。

可以说 M_0 对所有空间复杂度为 $S_1(n)$ 的 DTM 实行了对角线化, 因为可想象一个无穷的二维表, 其中第 ij 个元素表示第 i 台图灵机是否接受输入 j , 易见 M_0 的行为将在对角线上和某些图灵机不一致。特别地, 和那些在空间 $S_1(n)$ 内接受它们的输入的图灵机不一致。由引理 10.1 的推论 3, 存在和 M_0 等价的具有相同空间复杂度的单带 DTMM $_0'$ 。因为 M_0' 本身在表中 (即存在某个 k 值, M_0' 是第 k 台图灵机), 而 M_0' 不能和自己不一致, 因此可得, M_0 和 M_0' 的空间复杂度不为 $S_1(n)$ 。由于期望 M_0 的空间复杂度是 $S_2(n)$, 这里 $S_2(n)$ 和 $S_1(n)$ 是几乎相同的, 因此实际构造 M_0 是困难的任务。

定义 令 $f(n)$ 是一个任意函数。 $\inf_{n \rightarrow \infty} f(n)$ 是 n 趋近无穷时, $f(n), f(n+1), f(n+2), \dots$ 的最大下界。

例 11.1 由于 $(n^2 + 1)/n^2$ 是单调递减的,

$$\inf_{n \rightarrow \infty} \frac{n^2 + 1}{n^2} = \lim_{n \rightarrow \infty} \frac{n^2 + 1}{n^2} = 1$$

对于第 2 个例子, 如果 n 不是 2 的幂, 令 $f(n) = 1 - 1/n$; 如果 n 是 2 的幂, $f(n) = n$ 。则 $\inf_{n \rightarrow \infty} f(n) = 1$, 因为 $f(n), f(n+1), f(n+2), \dots$ 的最大下界是 $1 - 1/n$, 若 n 不是 2 的幂; 或是 $1 - 1/(n+1)$, 若 n 是 2 的幂。□

定理 11.1 令 $S_1(n) \geq n, S_2(n) \geq n$ 是两个空间可构造函数, 且

$$\inf_{n \rightarrow \infty} \frac{S_1(n)}{S_2(n)} = 0$$

则存在语言 L , 可被空间复杂度为 $S_2(n)$ 的 DTM 接受, 但不能被空间复杂度为 $S_1(n)$ 的 DTM 接受。[○]

○ 在文中, 读者通常会发现图灵机的空间复杂度的定义中忽略了输入带所扫描的单元数。然而, 不允许修改输入带上的符号, 在这种模型中, 讨论小于 n 的空间复杂度函数是有意义的, 此时条件 $S_1(n) \geq n$ 和 $S_2(n) \geq n$ 可以从定理的假设中移去。因为本书仅处理大空间复杂度, 所以任何小于 n 的复杂度不值得详细讨论, 因此就忽略了。这个定理可以放松 $S_1(n)$ 是空间可构造的这一限制。

证明: 令 M_0 是 4 带 DTM, 对长度为 n 的输入串 x 执行如下操作。

1. M_0 在每条带上标记出 $S_2(n)$ 个单元格。此后, 如果 M_0 的任意磁头企图离开标记的单元, 则 M_0 以不接受该输入串停机。

2. 如果 x 不是某个单带 DTM 的编码, 则 M_0 以不接受该输入串停机。

3. 否则, 令 M 是编码为 x 的 DTM。 M_0 确定了 M 使用的带符号的数目 t 和状态数 s 。 M_0 的第 3 条带可以作为计算 t 的“工作”内存。 M_0 在第 2 条带上写下 $S_1(n)$ 块, 每块包含 $\lceil \log t \rceil$ (单元, 各块之间用包含符号 # 的单元隔开, 即: 如果 $(1 + \lceil \log t \rceil) S_1(n) \leq S_2(n)$, 则共有 $(1 + \lceil \log t \rceil) S_1(n)$ 个单元。每个在 M 带子上出现的带符号都将编码为 M_0 第 2 条带上相应块中的二进制数。初始时, M_0 将其输入以二进制编码的形式放置在第 2 条带的块中, 对未使用的块则用空白符的编码进行填充。

4. 在带 3 上, M_0 建立一个有 $\lceil \log s \rceil + \lceil \log S_1(n) \rceil + \lceil \log t \rceil S_1(n)$ 单元的块, 块上的所有单元初始化为 0, 再次假定单元数不超过 $S_2(n)$ 。带 3 作为计数器计数至 $s S_1(n) t^{S_1(n)}$ 。

5. M_0 仿真 M , 使用带 1 即输入带, 来确定 M 的移动, 使用带 2 仿真 M 的带。使用带 3 的块以二进制对 M 的移动计数, 使用带 4 存储 M 的状态。如果 M 接受, 则 M_0 以不接受停机; 如果 M 以不接受停机, 或 M 的仿真试图使用多于带 2 分配的单元, 或带 3 的计数溢出, 即 M 的移动次数超出了 $s S_1(n) t^{S_1(n)}$, 则 M_0 接受。

上面描述的图灵机 M_0 以空间复杂度 $S_2(n)$ 接受某个语言 L 。假定 L 被某个 DTMM M_i 以空间复杂度 $S_1(n)$ 接受。由引理 10.1 的推论 3, 可假定 M_i 是单带图灵机。令 M_i 有 s 个状态和 t 个带符号。将 M_i 的 5 元组排成一行, 如果需要, 则加上前缀 1, 可以发现一个长为 n 的表示 M_i 的字符串 w , 这里 n 足够大, 使得 $S_2(n)$ 大于 $\text{MAX}[1 + \lceil \log t \rceil) S_1(n), \lceil \log s \rceil + \lceil \log S_1(n) \rceil + \lceil \log t \rceil S_1(n)]$ 。 $\inf_{n \rightarrow \infty} [S_1(n) / S_2(n)] = 0$ 的事实保证可以找到这样的 n 。接着, 给定 w 作为输入, M_0 有足够的空间仿真 M_i 。 M_0 接受 w 当且仅当 M_i 不接受 w 。但已经假定 M_i 接受 L , 即对所有输入, M_0 和 M_i 是一致的。由此可得, M_i 是不存在的, 即 L 不能被任何空间复杂度为 $S_1(n)$ 的 DTM 所接受。 \square

定理 11.1 的典型应用是用来说明, 例如, 存在一个语言可被空间复杂度 $n^2 \log n$ 的 DTM 所接受, 但不能被空间复杂度为 n^2 的 DTM 所接受。其他应用则见本章的剩余部分。

11.3 一个需要指数时间和空间的问题

10.6 节考虑了正则表达式的补的空性问题, 说明它是多项式空间完全的。由于正则集的类关于交和补运算是封闭的, 将交运算符 \cap 和补运算符 \neg 添加到正则表达式中并不会提升其所描述的集合的类。然而, 使用交和补运算符会大大缩短描述某些正则集的正则表达式的长度。由于具备浓缩的能力, 当时间的度量是以给定表达式的长度为基础时, 对于一些问题, 扩展正则表达式需要比“通常的”正则表达式更多的时间来求解。

定义 字母表 ϵ 的一个扩展正则表达式可定义如下:

1. ϵ, \emptyset, a (其中 a 属于 ϵ) 是表示为 $\{\epsilon\}, \{\emptyset\}$ 和 $\{a\}$ 的扩展正则表达式;
2. 如果 R_1 和 R_2 分别是表示语言 L_1 和 L_2 的扩展正则表达式, 则 $(R_1 + R_2), (R_1 \cdot R_2), (R_1^*), (R_1 \cap R_2)$ 和 $(\neg R_1)$ 也是扩展正则表达式, 分别表示语言 $L_1 \cup L_2, L_1 L_2, L_1^*, L_1 \cap L_2$ 和 $\Sigma^* - L_1$ 。

如果假定运算符的优先级按递增序为

$$+ \quad \cap \quad \neg \quad \cdot \quad *$$

可进一步将从扩展正则表达式中消去。

进一步, 运算符 \cdot 通常可以省略。例如 $a \neg b^* + c \cap d$ 意思是

$$((a \cdot (\neg(b^*))) + (c \cap d))$$

一个扩展的正则表达式如果不包含补符号则称为半扩展的。半扩展正则表达式的补的空性问题是确定给定的表达式 R 所表示的集合的补是否为空(等价的说法是, R 是否表示输入字母表上的所有字符串)。例如, 正则表达式

$$b^* + b^*aa^*(\epsilon + b(a+b)(a+b)^*) + b^*aa^*b$$

表示 $(a+b)^*$, 原因是

$$b^* + b^*aa^*(\epsilon + b(a+b)(a+b)^*) = \neg b^*aa^*b$$

现在开始证明半扩展正则表达式的补的空性问题不在 \mathcal{P} -SPACE 中, 由于 \mathcal{P} -SPACE 包含 \mathcal{NP} -TIME, 因此也不在 \mathcal{NP} -TIME 中。事实上, 我们将说明存在长度为 n 的半扩展正则表达式, 至少需要 $c'c^{\sqrt{n/\log n}}$ 空间(和时间)来求解补的空性问题, 这里常数 $c' > 0$, $c > 1$ 。

在 11.4 节, 我们可以看到完全扩展的正则表达式的补的空性问题比半扩展的正则表达式难得多。为了预期的结果, 本节将给出许多完全扩展的正则表达式的结果。

除了 10.6 节中使用的简写法外, 这里将使用 $\sum_{i=1}^k R_i$ 表示扩展正则表达式 $R_1 + R_2 + \cdots + R_k$, 还会使用 R^+ 替代 RR^* 。当谈及表达式的长度时, 是指由这些简写形式所表示的实际表达式的长度。

现在引入衡量尺度的记法。令 Σ 是一个字母表, x 是属于 Σ^* 的任意字符串。则 $\text{CYCLE}(x) = \{zy \mid x = yz, y \in \Sigma^*, z \in \Sigma^*\}$ 。令 $\#$ 是不属于 Σ 的一个特殊标记符号。则集合 $\text{CYCLE}(x\#)$ 称为是长度 $|x\#|$ 的衡量尺度。即, 衡量尺度是串 $x\#$ 的所有循环排列的集合。在不引起混淆的情况下, 有时也将 $x\#$ 本身看成是“衡量尺度”。

下面将使用衡量尺度度量在有效的图灵机计算中 ID 的长度。首先说明某些长的衡量尺度可以使用相应较短的半扩展正则表达式表示。

引理 11.1 对每个 $k \geq 1$, 存在长度大于 2^k 的衡量尺度, 可以用半扩展正则表达式 R 表示, 使 $|R| \leq ck^2$, c 是独立于 k 的常量。

证明: 令 $A = \{a_0, a_1, \dots, a_k\}$ 是 $k+1$ 个不同字符的字母表。令 $x_0 = a_0a_0$, $x_i = x_{i-1}a_i x_{i-1}a_i$ ($1 \leq i < k$)。由此

$$x_i = (\cdots((a_0^2 a_1)^2 a_2)^2 \cdots a_i)^2$$

x_0 的长度为 2, x_i 的长度大于 x_{i-1} 的长度的两倍。因此, x_{k-1} 的长度大于或等于 2^k , $k \geq 1$ 。

令 $\text{CYCLE}(x_{k-1}a_k)$ 为期望的衡量尺度。仍需要说明存在一个表示集合 $\text{CYCLE}(x_{k-1}a_k)$ 的半扩展正则表达式。

对于 $0 \leq i < k$, 令 A_i 表示 $(a_0 + a_1 + \cdots + a_i)$, $A - A_i$ 表示 $(a_{i+1} + a_{i+2} + \cdots + a_k)$ 。令

$$\begin{aligned} R_0 = & a_0 a_0 [(A - A_0)^+ a_0^2]^* (A - A_0)^+ \\ & + a_0 [(A - A_0)^+ a_0^2]^* (A - A_0)^+ a_0 \\ & + [(A - A_0)^+ a_0^2]^* (A - A_0)^+ a_0 a_0 (A - A_0)^+ \end{aligned}$$

和

$$\begin{aligned} R_i = & A_{i-1}^+ a_i A_{i-1}^+ a_i [(A - A_i)^+ (A_{i-1}^+ a_i)^2]^* (A - A_i)^+ A_{i-1}^+ \\ & + a_i A_{i-1}^+ a_i [(A - A_i)^+ (A_{i-1}^+ a_i)^2]^* (A - A_i)^+ A_{i-1}^+ \\ & + A_{i-1}^+ a_i [(A - A_i)^+ (A_{i-1}^+ a_i)^2]^* (A - A_i)^+ A_{i-1}^+ a_i A_{i-1}^+ \\ & + a_i [(A - A_i)^+ (A_{i-1}^+ a_i)^2]^* (A - A_i)^+ A_{i-1}^+ a_i A_{i-1}^+ \\ & + [(A - A_i)^+ (A_{i-1}^+ a_i)^2]^* (A - A_i)^+ A_{i-1}^+ a_i A_{i-1}^+ a_i (A - A_i)^+ \end{aligned}$$

其中 $1 \leq i \leq k-1$ 。

我们断言半扩展正则表达式

$$R = R_0 \cap R_1 \cap \cdots \cap R_{k-1} \cap A_{k-1}^+ a_k A_{k-1}^+$$

表示 $\text{CYCLE}(x_{k-1}a_k)$ 。证明, 通过对 i 归纳可知 $R_0 \cap R_1 \cap \cdots \cap R_i$ 字符串的形式为

$$y_2[(A-A_i)^+x_i]^*(A-A_i)^+y_1 \\ + [(A-A_i)^+x_i]^*(A-A_i)^+x_i(A-A_i)^+$$

这里 $y_1y_2 = x_i$ 。

归纳基础: 对于 $i=0$, R_0 精确表示了形式为

$$a_0^j[(A-A_0)^+a_0^2]^*(A-A_0)^+A_0^{2-j} \\ + [(A-A_0)^+a_0^2]^*(A-A_0)^+a_0a_0(A-A_0)^+$$

的字符串, $0 \leq j \leq 2$ 。因为 $x_0 = a_0$, 所以结果是显然的。

归纳步 假定 $R_0 \cap R_1 \cap \cdots \cap R_{i-1}$ 表示了所有形式为

$$y'_2[(A-A_{i-1})^+x_{i-1}]^*(A-A_{i-1})^+y'_1 \\ + [(A-A_{i-1})^+x_{i-1}]^*(A-A_{i-1})^+x_{i-1}(A-A_{i-1})^+ \quad (11-1)$$

的字符串, 这里 $y'_1y'_2 = x_{i-1}$ 。

考虑式(11-1)和 R_i 交集集中的字符串 z 。 A_{i-1}^+ 中的最长子串 z 必定是 x_{i-1} , 除非其位于串的头或尾, 在这种情况下, 是 y'_2 或 y'_1 。因此, R_i 中的 A_{i-1}^+ 可以替换为 x_{i-1} (除非其位于正则表达式的头或尾) 和 R_i 。位于 R_i 头或尾的 A_{i-1}^+ 和 A_{i-1}^+ 可以分别被替换为 y'_2 和 y'_1 。注意 $x_{i-1}a_ix_{i-1}a_i = x_i$, 可得

$$R_0 \cap R_1 \cap \cdots \cap R_i = y'_2a_ix_{i-1}a_i[(A-A_i)^+x_i]^*(A-A_i)^+y'_1 \\ + a_ix_{i-1}a_i[(A-A_i)^+x_i]^*(A-A_i)^+x_{i-1} \\ + y'_2a_i[(A-A_i)^+x_i]^*(A-A_i)^+x_{i-1}a_iy'_1 \quad (11-2) \\ + a_i[(A-A_i)^+x_i]^*(A-A_i)^+x_{i-1}a_ix_{i-1} \\ + [(A-A_i)^+x_i]^*(A-A_i)^+x_{i-1}a_ix_{i-1}a_i(A-A_i)^+$$

最后, 式(11-2)可以重写为

$$y_2[(A-A_i)^+x_i]^*(A-A_i)^+y_1 \\ + [(A-A_i)^+x_i]^*(A-A_i)^+x_i(A-A_i)^+$$

其中 $y_1y_2 = x_i$ 。可使用归纳假设。令 $i=k-1$ 并将上式与 $A_{k-1}^+a_kA_{k-1}^+$ 进行交操作, 由结果可知 R 表示 $\text{CYCLE}(x_{k-1}a_k)$ 。

剩下要证明的是 R 所表示的扩展正则表达式的长度界限是 ck^2 。该结论由以下事实立即可得: 存在常量 c_1 , 使得由每个缩写的 R_i 所表示的正则表达式的长度界限是 c_1k 。即, A_i 和 $A-A_i$ 表示长度为 $O(k)$ 的扩展正则表达式。因此, 每个定义 R_i 的项的和表示长度为 $O(k)$ 的扩展正则表达式。因此, 半扩展正则表达式 R 的长度界限是 ck^2 , 其中 c 是某个新的常量。 \square

使用类似引理 11.1 的构造过程, 可写出简写的正则表达式以表示除了 $x = x_{k-1}$ 外的所有 A^+ 中的串。该正则表达式使用刚构造的对应于 $\text{CYCLE}(x\#)$ 的半扩展正则表达式 R 。

引理 11.2 设字母表 $A = \{a_0, a_1, \dots, a_k\}$ 与字符串 x_i 如引理 11.1, 则存在一个表示 $\neg x_{k-1}$ 的长度为 $O(k^2)$ 的正则表达式。

证明: 串 x_{k-1} 遵循下述规则。

- i) 非空并以 a_0 开始。 a_0 成双出现, 并且每次成双出现之后跟随着一个 a_1 。
- ii) a_i 也是成双出现的, $1 \leq i \leq k-2$ 。每对出现的 a_i 都被 $a_0A_{i-1}^+$ 中的串分隔, 并且对中的第 2 个 a_i 之后跟随着 a_{i+1} 。即 a_i 的左边是 A_{i-1} 中的专有符号, 或直接在它左边的是 $A-A_i$ 中的符号接着是 a_0 后接 A_{i-1}^+ 中的字符串。

- iii) 恰好有两个 a_{k-1} 的实例, 第一个后接一个 a_0 , 第二个在串的右端。

更重要的是, x_{k-1} 是唯一的遵循这些规则的串。通过对 j 进行归纳可得, 如果 $b_1b_2 \cdots b_j$ 是满足

上述3条规则的一个串的前缀, 则 $b_1 b_2 \cdots b_j$ 是 x_{k-1} 的前缀。例如, 由规则(i), 第一个符号是 a_0 。同样由规则(ii), 单独的 a_0 不会终结一个字符串或后接一个不是 a_0 的符号, 因此字符串由 $a_0 a_0$ 开始。接着由规则(ii), $i=1$ 时, 串 $a_0 a_0 a_1$ 必定后接 a_0 , 依此类推。形式的归纳证明留给读者作为练习。由于 x_{k-1} 事实上满足规则, 可以看到其必定是满足这些规则的唯一串。

容易写出表示不满足上述3条规则之一或更多的串的正则表达式。不满足规则(i)的串表示为[○]

$$S_1 = \epsilon + (A - a_0)A^* + A^*a_0a_0[(A - a_1)A^* + \epsilon] + [\epsilon + A^*(A - a_0)]a_0[(A - a_0)A^* + \epsilon]$$

S_1 的第1项表示空串, 第2项表示不是从 a_0 开始的串, 第3项表示包含一对 a_0 但其后不跟随 a_1 的串, 最后一项是包含单独的 a_0 的串。那些不满足规则(ii)的串则由下式表示

$$S_2 = \sum_{i=1}^{k-2} [A^*a_iA_{i-1}^*a_i[(A - a_{i+1})A^* + \epsilon] + [\epsilon + A^*(A - A_i)]A_{i-1}^*a_i[(A - a_0)A^* + \epsilon]$$

最后, 不满足规则(iii)的串是

$$S_3 = A^*a_{k-1}A^*a_{k-1}A^* + (A - a_{k-1})^*a_{k-1}(A - a_{k-1})^* + (A - a_{k-1})^* + (A - a_{k-1})^*a_{k-1}(A - a_0)A^*$$

在 S_3 中, 第1项除了包括正好含有两个 a_{k-1} 的串外, 还包括含有多于两个 a_{k-1} 的所有串, 其中第二个 a_{k-1} 不位于串的右端。因此 $S_1 \cup S_2 \cup S_3$ 表示 $\neg x_{k-1}$ 。易见 $S_1 \cup S_2 \cup S_3$ 的长度是 $O(k^2)$ 。 \square

在说明半扩展正则表达式补的空性问题不存在多项式空间界限或多项式时间界限的算法之前, 先进行两个基本的考察。

定义 从 Σ_1^* 到 Σ_2^* 的同态映射 h 是一个函数, 对任意串 x 和 y , 有 $h(xy) = h(x)h(y)$ 。由此, $h(\epsilon) = \epsilon$, $h(a_1 a_2 \cdots a_n) = h(a_1)h(a_2) \cdots h(a_n)$ 。因此, 同态映射 h 由 Σ_1 中的每个 a 的 $h(a)$ 值所唯一定义。

若对于 Σ_1 中的每个 a , $h(a)$ 的值是 Σ_2 中的单一符号, 则称同态映射 h 是保持长度的。保持长度的同态映射仅仅是对符号重命名, 可通过重命名将几个符号看成同一符号。如果 w 在 Σ_2^* 中, 则 $h^{-1}(w) = \{x \mid h(x) = w\}$ 。如果 $L \subseteq \Sigma_2^*$, 则 $h^{-1}(L) = \{x \mid h(x) \in L\}$ 。

例 11.2 令 $\Sigma_1 = \{a, b, c\}$, $\Sigma_2 = \{0, 1\}$ 。如下定义同态映射 h , $h(a) = 010$, $h(b) = 1$, $h(c) = \epsilon$, 则 $h(abc) = 0101$ 。可见 h 不是保持长度的同态映射。令 L 是由扩展正则表达式 $1 + \neg 1^*$ 或由等价的一般正则表达式 $1 + (0+1)^*0(0+1)^*$ 表示的语言, 则 $h^{-1}(L)$ 可由扩展正则表达式 $c^*bc^* + \neg(b+c)^*$ 表示, 或由一般正则表达式 $c^*bc^* + (a+b+c)^*a(a+b+c)^*$ 表示。 \square

引理 11.3 令 h 是保持长度的从 Σ_1^* 到 Σ_2^* 的同态映射, R_2 是表示集合 $S \subseteq \Sigma_2^*$ 的一个扩展正则表达式, 则可以构造一个表示 $h^{-1}(S)$ 的扩展正则表达式 R_1 , R_1 的长度界限是 R_2 的长度的常数倍(仅依赖于 h), 并且仅当 R_2 包含一个补符号时 R_1 也包含一个补符号。

证明: 将 R_2 中 Σ_2 的每个符号出现替换为一个正则表达式, 表示所有由 h 映射为该符号的集合。例如, 如果 $\{a_1, a_2, \cdots, a_r\}$ 是映射为 a 的符号集合, 则将 a 替换为 $(a_1 + a_2 + \cdots + a_r)$ 。令 R_1 是得到的扩展正则表达式。 R_1 表示 $h^{-1}(S)$ 的证明则可通过对 R_2 中运算符 $+$, \cdot , $*$, \cap 和 \neg 的出现数目进行归纳而得。 \square

现在必须如同引理 10.2 那样讨论图灵机的瞬时描述序列的表示。即给定图灵机 $M = \{Q, T, I, \delta, b, q_0, q_f\}$, 通过 T 中的符号序列加上一个形为 $[qX]$ 的其他符号来表示一个 ID, 这里 $q \in Q$, $x \in T$ 分别表示状态和输入磁头的位置。如果需要, 可使用一些空格来拉长一个 ID 以使同一个计算中的所有 ID 有相同的长度。

○ 使用 $A - a_i$ 表示 $\sum_{j=0, j \neq i}^k a_j$ 。

为了比较一个 ID 中的符号和后继 ID 中的“相应”符号,我们将一个表示 ID 长度的衡量尺度与 ID 本身放置在一起。形式化地说,使用一个“双轨”的符号串,上轨是 ID 序列,下轨则包含重复的衡量尺度。即“双轨”符号是对 $[a, b]$, a 是上轨中的符号, b 是下轨中的符号。该格局如图 11-1 所示,其中使用的衡量尺度是 $CYCLE(x\#)$ 。(x 是不包含 $\#$ 的任意字符串。)

定义 给定 TM M , 令 $\Delta_1 = T \cup (Q \times T) \cup \{\#\}$, Δ_2 是在 $x\#$ 中出现的符号集合。即 Δ_1 是可以出现在上轨中的符号集合, Δ_2 是可以出现在下轨中的符号集合。“双轨”的字母表是 $\Delta_1 \times \Delta_2$ 。衡量尺度为 $CYCLE(x\#)$ 的 M 的一个有效计算是如图 11-1 所示的 $(\Delta_1 \times \Delta_2)^*$ 中的串, M 的一个移动可使 $C_i \vdash C_{i+1}$, $0 \leq i < f$, C_0 是初始 ID, 在最后一个 ID C_f 中状态是 q_f 。在 ID 的末尾填充空白使所有 ID 的长度皆为 $|x|$ 。

上轨	#	C_0	#	C_1	#	...	#	C_f	#
下轨	#	x	#	x	#	...	#	x	#

图 11-1 带衡量尺度的 ID 序列

引理 11.4 令 R_1 是相应于 $CYCLE(x\#)$ 的扩展正则表达式, R'_1 是 $x\#$ 的字母表上表示除了 x 外所有串的正则表达式。可以构造一个表示 M 的衡量尺度为 $CYCLE(x\#)$ 的一个无效计算的扩展正则表达式 R_2 , R_2 的长度与 $|R_1| + |R'_1|$ 成线性比例, 比例常数仅依赖于 M 。仅当 R_1 或 R'_1 包含一个补符号时, R_2 包含一个补符号。

证明: 一个串是衡量尺度为 $CYCLE(x\#)$ 的一个无效计算, 当且仅当

1. 下轨不在 $\#(x\#)^*$ 中, 或者
2. 下轨在 $\#(x\#)^*$ 中, 但上轨不是一个有效的计算。

如果 Δ_1 是上轨的字母表, Δ_2 是下轨的字母表, 令 h_1 和 h_2 分别是 $(\Delta_1 \times \Delta_2)^*$ 到 Δ_1^* 和到 Δ_2^* 的同态映射, 满足 $h_1([a, b]) = a$, $h_2([a, b]) = b$ 。则

$$h_2^{-1}[\Delta_2^* \# (R'_1 \cap (\Delta_2 - \#)^*) \# \Delta_2^* + (\Delta_2 - \#) \Delta_2^* + \Delta_2^* (\Delta_2 - \#)] + \epsilon$$

仅表示所有下轨不在 $\#(x\#)^*$ 中的串。 h_2^{-1} 参数的第一项表示那些两个 $\#$ 之间不是 x 而是其他东西的串, 其他项则表示不是以 $\#$ 开始或结束的串。由引理 11.3, 存在一个长度与表示该集合 $|R'_1|$ 成比例的扩展正则表达式。

一个给定的串满足条件 2 是因为由等于衡量尺度的若干符号所分隔的两个符号并不反映 M 的移动, 或因为发生了一个或若干个下面的格式错误。

- i) 在上轨中, 串不是以 $\#$ 开始或结束;
- ii) 在第一个 ID 中, 没有状态或有两个或更多的状态出现;
- iii) 第一个 ID 没有包含初始状态作为其第一个符号的成分;
- iv) 接受状态不作为任何符号的成分出现;

v) 第一个 ID 的长度不正确。即上轨中的前两个 $\#$ 不像下轨中的前两个 $\#$ 一样位于相同的单元格中。

下面将看到对于 $(\Delta_1 \times \Delta_2)^*$ 中不能反映 M 的合法移动的序列如何写扩展正则表达式。如同引理 10.2, 我们注意到在有效的计算中, 上轨中的 3 个连续符号 $c_1 c_2 c_3$ 唯一确定放置于符号 c_2 右边的符号 $|x\#|$ 。与以前一样, 设这个符号为 $f(c_1 c_2 c_3)$ 。令

$$R_{c_1 c_2 c_3} = (\Delta_1 + \Delta_2)^* [h_1^{-1}(c_1 c_2 c_3 \Delta_1^*) \cap h_2^{-1}(R_1)] [h_1^{-1}(\Delta_1 (\Delta_1 - f(c_1 c_2 c_3)) \Delta_1)] (\Delta_1 \times \Delta_2)^* \quad (11-3)$$

再令

$$R = \sum_{c_1 c_2 c_3} R_{c_1 c_2 c_3}$$

注意, $h_1^{-1}(c_1 c_2 c_3 \Delta_1^*)$ 表示 $(\Delta_1 \times \Delta_2)^*$ 中的所有上轨以 $c_1 c_2 c_3$ 开始的串。 $h_2^{-1}(R_1)$ 表示下轨是正确的且长为 $|x\#|$ 的串。由此, 它们的交是所有长为 $|x\#|$ 的串, 其下轨包含 $\text{CYCLE}(x\#)$ 中的串, 且上轨以 $c_1 c_2 c_3$ 开始。基于这些观察, 很清楚, R 包括了所有满足条件 2 的串, 它们反映了 M 在某点进行了非法的移动。 R 也包括下轨不在 $\#(x\#)^*$ 中的串; 这些串满足条件 1, 它们在 R 中出现或不出现不是实质的。进一步, R 的长度是 R_1 长度的常数倍(依赖于 M)。为了说明其真实性, 注意到 h_2^{-1} 将正则表达式扩展了与 M 相关的常数倍。另外, h_1^{-1} 最多将表达式扩展了 $3 \parallel \Delta_2 \parallel$ 倍, 原因是对于恰好 k 个 b 的值, 如果 $h_1(b) = a$, 则 $h_1^{-1}(a) = (b_1 + b_2 + \dots + b_k)$ 的长度为 $2k+1$ 。由于 $1 \leq k \leq \parallel \Delta_2 \parallel$ 是显然的, 所以 $|h_1^{-1}(a)| \leq 3 \parallel \Delta_2 \parallel$ 。更进一步, 因为 Δ_2 中的每个符号都在 R_1 中出现, 很清楚有 $\parallel \Delta_2 \parallel \leq |R_1|$ 。因此, 式(11-3)中的项 $h_1^{-1}(c_1 c_2 c_3 \Delta_1^*)$ 和 $h_1^{-1}(\Delta_1 (\Delta_1 - f(c_1 c_2 c_3)) \Delta_1)$ 的长度是 $O(|R_1|)$ 。最后, 相应于 $h_2^{-1}(R_1)$ 和 $(\Delta_1 \times \Delta_2)^*$ 的项的长度很明显是 $O(|R_1|)$, 其中的常数仅依赖于 M 。因此, 式(11-3)的长度即 R 的长度是 $O(|R_1| + |R'_1|)$ 。

容易写出对应于格式错误的正则表达式, 这留给读者完成。按这种方式, 可以构造一个扩展正则表达式 R_2 , 其长度和 $|R_1| + |R'_1|$ 成线性比例, 且仅当 R_1 或 R'_1 包含一个补符号时 R_2 包含一个补符号。 \square

定理 11.2 确定半扩展正则表达式是否表示了其字母表上所有串的任意算法^①的空间复杂度(和时间复杂度)至少是 $c'c^{\sqrt{n \log n}}$, 其中 $c' > 0$ 和 $c > 1$ 是常数, n 的数目是无穷的。

证明: 令 L 是空间复杂度为 2^n 但不为 $2^n/n$ 的任意语言^②(由定理 11.1 可知这样的语言是存在的)。令 M 是接受语言 L 的 DTM。

假定有一个空间复杂度为 $f(n)$ 的 DTM M_0 , 可判别由一个半扩展正则表达式表示的集合的补是否空。则可以按照如下方式使用 M_0 识别语言 L 。令 $w = a_1 a_2 \dots a_n$ 是长为 n 的输入串。

1. 构造长度的衡量尺度为 $2^n + 1$ 或更大的半扩展正则表达式 R_1 。由引理 11.1(取 $k = n$), 存在一个长度为 $O(n^2)$ 的表达式 R_1 。进一步, R_1 可以只使用空间复杂度 $O(n^2)$ 构造。类似地, 构造表示 $\neg x$ 的 R'_1 , 这里 $R_1 = \text{CYCLE}(x\#)$ 。由引理 11.2, R'_1 的长度为 $O(n^2)$, 容易看出 R'_1 可使用空间 $O(n^2)$ 构造。

2. 构造一个半扩展正则表达式 R_2 以表示衡量尺度为 R_1 的 M 的无效计算。由引理 11.4, 存在一个长度至多为 $c_1 n^2$ 的 R_2 , c_1 为仅依赖于 M 的常量。

3. 构造一个正则表达式 R_3 以表示所有 $(\Delta_1 \times \Delta_2)^*$ 中的上轨不以 $\#[q_0 a_1] a_2 \dots a_n b \dots b\#$ 开始的所有串, 其中 q_0 是 M 的初始状态。很明显长度为 $O(n)$ 的表达式 R_3 是存在的。因此, $|R_2 + R_3| \leq c_2 n^2$, c_2 为常量。

4. 应用 M_0 于 $R_2 + R_3$ 以判别 $R_2 + R_3$ 的补是否空。如果不空, 则存在 M 的一个输入为 w 的有效计算, 因此 w 在 L 中。否则, w 不在 L 中。

可以构造一个空间复杂度为 $f(c_2 n^2 \log n)$ 的 DTM M' 以实现这个识别 L 的算法。 f 参数中的 $\log n$ 因子是由于正则表达式 $R_2 + R_3$ 是 n 个符号的字母表上的, 必须编码为固定的字母表。因为我们假设 L 是空间复杂度为 2^n 但不为 $2^n/n$ 的语言, 所以至少对某些 n 值必有 $f(c_2 n^2 \log n) > 2^n/n$ 。但对有限数目的 n , 若 $f(c_2 n^2 \log n)$ 超过了 $2^n/n$, 则存在上述识别算法的一个修正版本, 它首

① 假定使用类似于引理 10.3 中的编码方法。

② 2^n 的选择不是本质性的。可以将 2^n 替换为任意的指数函数 $f(n)$, $2^n/n$ 替换为任意的增长速度较 $f(n)$ 慢的函数, 只要每个都是空间可构造的即可。

先通过有限的“表查看”检查 $|w|$ 是否使 $f(c_2 n^2 \log n) > 2^n/n$ 成立的一个 n 值。如果是, w 是否在 L 中。因此, 对于无穷数目的 n , $f(c_2 n^2 \log n)$ 超过了 $2^n/n$, 这样, 对于无穷数目的 m 和常量 $c_3 > 0$, $c > 0$ 和 $c' > 1$, $f(m) \geq 2^{c_3 \sqrt{m/\log m}} / \sqrt{m/\log m} \geq c(c')^{\sqrt{m/\log m}}$ 。□

推论 半扩展正则表达式的等价问题需要 $c'c^{\sqrt{n/\log n}}$ 时间和空间开销, 其中常数 $c' > 0$, $c > 1$, n 的数目是无穷的。

证明: 因为表示所有串的正则表达式是简短且容易写出的, 容易证明补集为空的问题可多项式归约为等价问题。□

11.4 一个非基本的问题

现在考虑整个扩展正则表达式类。因为有补算子, 仅需要考虑空问题, 即给定的扩展正则表达式 R 是否表示一个空集。可以看到使用补算子可以更加紧凑地表示正则集, 另外扩展正则表达式的空性问题比半扩展正则表达式的补集为空的问题更难。

定义函数 $g(m, n)$ 如下:

$$1. g(0, n) = n,$$

$$2. g(m, n) = 2^{g(m-1, n)}, m > 0,$$

由此 $g(1, n) = 2^n$, $g(2, n) = 2^{2^n}$, 并且

$$g(m, n) = 2^{2^{\dots^{2^n}}}$$

式中有 m 个 2, 最后是 2 的 n 次方。称函数 $f(n)$ 是基本的 (elementary), 如果对于所有 n (除了一个有限集合外) 其上界是 $g(m_0, n)$, 其中 m_0 是某个固定的整数。

11.3 节的技术可用来说明不存在基本函数 $S(n)$ 使整个扩展正则表达式类的空性问题的空间复杂度为 $S(n)$ 。在处理之前, 对衡量尺度为 $\text{CYCLE}(x\#)$ 的图灵机 $M = (Q, T, I, \delta, b, q_0, q_t)$ 的有效计算的定义稍做修改。删去第一个标识符 $\#$, 将最后一个标识符改为一个新的符号 $\$$, 如图 11-2 所示。 C 是 ID (如果需要, 则拉伸与 x 等长)。 M 的一个移动可使 $C_i \vdash C_{i+1}$, $0 \leq i < f$, C_0 是初始 ID, 在最后一个 C_i 中状态是 q_t 。

使用下面的惯例记法。假定 $x \in \Sigma^*$ 。

1. $\Delta_1 = T \cup (Q \times T) \cup \{\#, \$\}$ 是上轨字母表;
2. $\Delta_2 = \Sigma \cup \{\#, \$\}$ 是下轨字母表;
3. $h_1: \Delta_1 \times \Delta_2 \rightarrow \Delta_1$, 其中 $h_1([a, b]) = a$;
4. $h_2: \Delta_1 \times \Delta_2 \rightarrow \Delta_2$, 其中 $h_2([a, b]) = b$ 。

上轨	C_0	$\#$	C_1	$\#$	\dots	$\#$	C_f	$\$$
下轨	x	$\#$	x	$\#$	\dots	$\#$	x	$\$$

图 11-2 带衡量尺度的有效计算的新格式

引理 11.5 令 R_1 是一个对应于 $\text{CYCLE}(x\#)$ 的扩展正则表达式。可以构造一个扩展正则表达式 R_2 , 表示衡量尺度为 $\text{CYCLE}(x\#)$ 的图灵机 M 的有效计算的循环排列的集合, 满足 $|R_2|$ 等于 $O(|R_1|)$, 常数因子仅依赖于 M 。

证明: 一个串是衡量尺度为 $\text{CYCLE}(x\#)$ 的有效计算的循环排列, 当且仅当

1. 下轨是形为 $(x\#)^* x \$$ 的串的一个循环排列;
2. 上轨是 M 的有效计算的一个循环排列;

3. 上轨和下轨循环排列的次数是一样的。

令 R'_1 是 R_1 中的 $\#$ 替换为 $\$$ 所得的扩展正则表达式。满足条件 1 的串可用 $h_2^{-1}(U_1 \cap U_2 \cap U_3)$ 表示, 这里

$$U_1 = \Sigma^* (\# + \$) [(R_1 + R'_1) \cap (\Sigma^* \# + \Sigma^* \$)]^* \Sigma^*$$

$$U_2 = (\Sigma + \#)^* \$ (\Sigma + \#)^*$$

$$U_3 = (R_1 + R'_1)^*$$

U_1 中的子表达式 $[(R_1 + R'_1) \cap (\Sigma^* \# + \Sigma^* \$)]$ 表示两个串 $x\#$ 和 $x\$$ 。因此, U_1 表示 $\Sigma^* (\# + \$) (x\# + x\$)^* \Sigma^*$ 。表达式 U_2 仅允许包含一个 $\$$ 的串。因此, $U_1 \cap U_2$ 中的任意串有形式

$$y_2 \# x \# x \# \cdots x \# x \$ x \# \cdots x \# y_1$$

其中 y_1 和 y_2 属于 Σ^* 。表达式 U_3 使得 $|y_2| + |y_1| = |x|$, 由此容易看出, 对于 $U_1 \cap U_2 \cap U_3$ 中的串, 有 $y_1 y_2 = x$ 。因此, $S_1 = U_1 \cap U_2 \cap U_3$ 表示满足条件 1 的所有串的下轨。

对于条件 2, 引理 11.4 说明如何写对应于衡量尺度为 $\text{CYCLE}(x\#)$ 的 M 的有效计算的半扩展正则表达式。这些技术容易表示为图 11-2 所示的格式, 而表示下轨上衡量尺度正确的图灵机 M 的无效计算的循环排列的表达式 E 的工作则留给读者自行构造。对 E 应用补算子的结果是正则表达式 S_2 。这是唯一使用补算子之处。很清楚, $S_1 \cap S_2$ 表示的所有串都满足条件 1 和 2。

当条件 1 和 2 满足时, 相应于条件 3 的表达式是容易得到的。仅需检查在双轨中一个符号是 $\$$ 即可。将该表达式和 $S_1 \cap S_2$ 相交即可完成定理的证明。□

现在说明基于引理 11.5 的构造方法, 对于长度的稍许增长, 该如何表示越来越长的衡量尺度。

直观上说, 先对一个较小衡量尺度比如 n , 构造一个正则表达式, 接着使用该衡量尺度对所有正好接受长度为 n 的输入串的某个图灵机上的关于该衡量尺度的有效计算的循环排列构造一个新的正则表达式。仔细选择该图灵机, 使其对于长度为 n 的输入, 至少移动 2^{n+2} 次, 这样其有效计算的循环排列的集合本身是长为 2^n 或更大的衡量尺度。使用新的衡量尺度重复上述过程, 得到的是越来越长的衡量尺度, 至少是 $2^n, 2^{2^n}$ 等等。

令 M_0 是行为如下的一台特殊的单带图灵机。

1. M_0 通过扫描它的输入带直到遇到第一个空白格以检查它的输入带是否以一串 a 开始;
2. 如果其输入带是以 m 个 a 后随一个空白格的字符串开始, 则 M_0 对输入带上的 a 和随后的空白格部分以二进制形式从 0 到 $2^{m+1} - 1$ 计数;
3. 当 M_0 计数至 $2^{m+1} - 1$ 时, 停机接受。

注意下面关于 M_0 的事实。对每个 n , M_0 恰好接受的是长度为 n 的输入串, 即 a^n 。其特殊的有效计算包括至少执行 2^{n+2} 次移动, 一半用于位加法, 一半用于进位。最后, 当给定长度为 n 的输入时, M_0 仅扫描了 $n+1$ 个带单元。

所以, 考虑图 11-2 所示的衡量尺度为 $\text{CYCLE}(x\#)$ 的 M_0 的有效计算, 如果 $|x| = 1$, 则存在唯一的上轨以 $[q_0 a] a \cdots ab \#$ 起始的有效计算。由引理 11.5, 可以构造一个正则表达式 R_2 以表示衡量尺度为 $\text{CYCLE}(x\#)$ 的 M_0 的一个有效计算的循环排列。表达式 R_2 表示由 $(\Delta_1 \times \Delta_2)^*$ 中的固定串 w 的循环排列所组成的集合。串 $h_1(w)$ 是 M_0 针对输入 a^n 的有效计算, 这里的 $|x|$ 值, 回想一下, 现在是 $n+1$ 。因为对于输入 a^n , M_0 至少移动 2^{n+2} 次, 可以使用 R_2 本身作为长度至少为 2^{n+2} 的衡量尺度。

总之, 假设给定衡量尺度 $R_1 = \text{CYCLE}(x\#)$, 其中 $x \in \Sigma^*$ 。可以基于 R_1 构造一个衡量尺度至少为 $2^{|x|+1}$ 的特殊的图灵机。由引理 11.5, 对于该衡量尺度存在长度至多为 $c_1 |R_1|$ 的扩展正则表达式 R_2 , 这里 c_1 仅依赖于 M_0 。

引理 11.6 对所有的 $i \geq 1$ 和 $m \geq 2$, 存在长度至少为 $g(i, m)^\ominus$ 的衡量尺度, 可以用长度至多为 $c(c_1)^{i-1}m^2$ 的扩展正则表达式所表示, 这里 c_1 是上一段引入的依赖于 M_0 的常数, c 是引理 11.1 中的常数。

证明: 通过对 i 归纳, 可知对某个 x 能找到一个表示 $\text{CYCLE}(x\#)$ 的扩展正则表达式 R_1 , 这里 $|x\#| \geq g(i, m)$ 且 $|R_1| \leq c(c_1)^{i-1}m^2$ 。归纳基础 $i=1$, 由引理 11.1 显然成立。可构造出表示 $\text{CYCLE}(x\#)$ 的长度为 cm^2 的扩展正则表达式。

对于归纳步, 假定有表示 $\text{CYCLE}(x\#)$ 的扩展正则表达式 R_1 , $|R_1| \leq c(c_1)^{i-2}m^2$ 且 $|x\#| \geq g(i-1, m)$ 。由上面关于 DTM M_0 的分析讨论, 可以构造长度至多为 $c_1|R_1| = c(c_1)^{i-1}m^2$ 的扩展正则表达式 R 。 R 表示 $\text{CYCLE}(y\$)$, 这里 $|y\$| \geq 2^{|x\#|} \geq 2^{g(i-1, m)} = g(i, m)$ 。 \square

定理 11.3 令 $S(n)$ 是任意的基本函数, 则不存在 $S(n)$ 空间界限的 [因此也没有 $S(n)$ 时间界限的] DTM 以判别扩展正则表达式是否表示空集。

证明: 使用反证法。假定空间界限为 $g(k_0, n)$ 的 TMM_1 可以判别一个扩展正则表达式是否表示空集。由定理 11.1, 存在语言 L 可被空间界限为 $g(k_0+1, n)$ 的 DTMM 所接受, 但不能被空间界限为 $[g(k_0+1, n)/n]$ 的 DTM 所接受。由假设, M_1 存在, 因此可以设计一个行为如下的 DTMM_2 以识别语言 L 。

1. 给定长度为 n 的输入串 $w = a_1a_2 \cdots a_n$, M_2 构造一个长度为 d_1n^2 , 可表示至少衡量尺度为 $g(k_0+1, n)$ 的扩展正则表达式 R_1 , 这里 d_1 是独立于 n 的常量。可使用引理 11.6 提示的算法来构造 R_1 。

2. 接着, M_2 构造一个扩展正则表达式 R_2 以表示接受语言 L 的 $g(k_0+1, n)$ 空间界限的 TMM 的有效计算的集合。由引理 11.5, 可使 $|R_2| \leq d_2|R_1|$, 其中 d_2 是常数。

3. 随后, M_2 构造扩展正则表达式 R_3 以表示衡量尺度为 R_1 , 初始 ID $C_0 = [q_0a_1]a_2 \cdots a_1bb \cdots$ 的 M 的有效计算, 其中 q_0 是 M 的初始状态。即

$$R_3 = R_2 \cap h_1^{-1}([q_0a_1]a_2 \cdots a_nb^*\#)(\Delta_1 \times \Delta_2)^*$$

这里 h_1 是引理 11.5 中的同态映射, $\Delta_1 \times \Delta_2$ 是 R_2 的字母表。由引理 11.3, $|R_3| \leq |R_2| + d_3n$, 其中 $d_3 > 0$ 是常数, 因此

$$|R_3| \leq d_1d_2n^2 + d_3n \quad (11-4)$$

4. 最后, M_2 将 R_3 编码为如定理 11.2 中所示的固定字母表, 并使用 M_1 测试扩展正则表达式 R_3 是否表示空集。如果答案为是, 则 M_2 拒绝 w ; 如果答案为否, 则 M_2 接受 w 。由此 M_2 接受语言 L 。

现在不难看出第 4 步是最大的空间开销步, M_1 需要空间 $g(k_0, |R_3| \log |R_3|)$ 以确定 R_3 是否表示空集。因此, M_2 也需要这样大小的空间。使用式 (11-4) 中 R_3 的界限, 可以看出 M_2 的空间复杂度是 $S(n) = d_4g(k_0, d_3n^2 \log n)$, 其中 d_4 和 d_5 是常数, 原因是对除有穷个例外的所有 n , 式 (11-4) 中的第 1 项即 $d_1d_2n^2$ 大于第 2 项 d_3n 。然而, 由定理 11.2 的讨论, 可以说明 M_2 的空间复杂度对于无穷个 n , 其值都超过 $g(k_0+1, n)/n$ 。因此, 如果 M_2 存在, 则有

$$g(k_0+1, n)/n < d_4g(k_0, d_3n^2 \log n) \quad (11-5)$$

对于无穷多个 n 值。但不管 d_4 和 d_5 取何值, 只有有穷个 n , 使式 (11-5) 成立, 对此, 读者应容易得出该结论。由此可得 M_2 是不存在的, 因此, M_1 也是不存在的。因此, 扩展正则表达式的空性问题对于任意的基本空间界限的图灵机是不可判定的。 \square

习题

11.1 称函数 $T(n)$ 是时间可构造的, 如果存在一个 DTM M , 对于给定的长度为 n 的输入, 其在

\ominus g 在 11.4 节定义。

停机前恰好移动了 $T(n)$ 次。说明下面的函数是时间可构造的。

a) n^2

b) 2^n

c) $n!$

*11.2 说明每个时间可构造函数是空间可构造的。如果 $S(n)$ 是空间可构造的, 则 $c^{S(n)}$ 是时间可构造的, 其中 c 是常整数。

*11.3 说明如果 L 可被一台 k 带 DTM (NDTM) 在时间 $T(n)$ 内接受, 则存在一台在时间 $O(T^2(n))$ 内接受 L 的单带 DTM (NDTM)。

*11.4 (DTM 的时间层次性。)说明如果 $T_1(n)$ 和 $T_2(n)$ 是时间可构造函数, 并且

$$\inf_{n \rightarrow \infty} \frac{T_1^2(n)}{T_2(n)} = 0$$

则 DTM 对某个语言, 可在时间 $T_2(n)$ 内接受该语言但不能在时间 $T_1(n)$ 内接受该语言。

[提示: 由习题 11.3, 只需要对时间复杂度为 $[T_1(n)]^2$ 的单带 DTM 进行对角线化即可, 对角线化可在一台多带 DTM 上进行。]

下面两个习题是关于弱形式的非确定图灵机层次性的。

**11.5 证明对于每个整数 $k \geq 1$, 存在可被空间复杂度为 n^{k+1} 的 NDTM 接受但不能被空间复杂度为 n^k 的 NDTM 接受的语言。

**11.6 对于时间复杂度, 证明和习题 11.5 相同的结果。

下面两个习题考虑 DTM 的时间复杂度的紧致性。

**11.7 证明每个可被时间复杂度为 $T(n)$ 的 k 带 DTM 所接受的语言也可被时间复杂度为 $O(T(n) \log T(n))$ 的 2 带 DTM 所接受。

11.8 使用习题 11.7 的结果说明如果 $T_1(n)$ 和 $T_2(n)$ 是时间可构造函数, 并且

$$\inf_{n \rightarrow \infty} \frac{T_1(n) \log T_1(n)}{T_2(n)} = 0$$

则 DTM 对某个语言, 可在时间 $T_2(n)$ 内接受该语言但不能在时间 $T_1(n)$ 内接受该语言。

*11.9 说明如果 L 可被一个空间复杂度为 $S(n)$ 和时间复杂度为 $T(n)$ 的 DTM (NDTM) M 所接受, c 是大于 0 的任意常数, 则 L 可被一个空间复杂度为 $\text{MAX}(cS(n), n+1)$ 和时间复杂度为 $\text{MAX}(cT(n), 2n)$ 的 DTM (NDTM) M' 接受。[提示: M' 必须先将 M 的单元块浓缩为其自身带上的单一单元。]

11.10 说明如果 L 可被一个时间复杂度为 $T(n)$ 的 NDTM 所接受, 则存在一个常数 c , 使得 L 可被一个时间复杂度为 $c^{T(n)}$ 的 DTM 所接受。

*11.11 设 T_1 和 T_2 是两个函数, 并且

$$\inf_{n \rightarrow \infty} \frac{T_1(n)}{T_2(n)} = 0$$

则 RAM 对某个语言, 可在时间 $T_2(n)$ 内但不能在时间 $T_1(n)$ 内接受该语言, 这里采用对数代价模型。

11.12 完成引理 11.2 的证明。

11.13 给定衡量尺度为 $\text{CYCLE}(x\#)$ 的扩展正则表达式 R_1 , 构造一个衡量尺度为 $\text{CYCLE}((x\#)^)$ 的扩展正则表达式 R_2 。 R_2 的长度的界限必须是 R_1 的长度的常数倍。

**11.14 给定一个 $O(n)$ 空间界限的判别一个半扩展正则表达式是否表示一个非空集的非确定算

法。如果尝试判别一个正则表达式是否表示所有串, 会失败吗? 为何必定失败?

- 11.15 给出一个指数时间复杂度的确定型算法以解决半扩展正则表达式的补集的空性问题。
- 11.16 写出引理 11.4 中“格式错误”(format errors)的正则表达式。
- 11.17 函数 $F(n)$ 的定义为: $F(0) = 1$, $F(n) = 2^{F(n-1)}$, $n \geq 1$ 。 $F(n)$ 是基本函数吗? (该函数在 4.7 节中首次引入。)
- *11.18 给定一个判别扩展正则集所表示的集合是否为空的算法。算法的时间和空间复杂度为何?
- **11.19 说明仅使用算子 $+$ 、 \cdot 和 \neg 的扩展正则表达式的空性问题不是基本的。

研究性问题

- 11.20 本章所提示的自然研究领域是寻找有实际意义的、可以证明是难的问题。这个方向包括 Fischer 和 Rabin[1974]关于基本算术复杂度的研究、Cook 和 Reckhow[1974]关于定理证明的研究以及 Hunt[1974]对于语言理论问题的讨论, 还有文献和注释中提及的其他作者的工作。
- 11.21 另外一个有意思的问题是对于 DTM 的时间层次性以及相应于 NDTM 的时间和空间层次性实际上是否比习题 11.5、11.6 和 11.8 所揭示的更严格。例如, 是否存在时间可构造函数 $T(n)$, 使得没有语言可在 $T(n)\log T(n)$ 内被识别但不能在 $T(n)$ 内被识别。读者可查阅 Seiferas、Fischer 和 Meyer[1973]以了解 NDTM 的已知最紧的层次性。

文献和注释

图灵机的时间和空间层次性的广泛研究可见 Hartmanis 和 Stearns[1965], Hartmanis、Lewis 和 Stearns[1965]。Rabin[1963]是关于时间复杂度的早期文章, 值得学习。习题 11.7 和 11.8 给出的时间层次性的改进取自 Hennie 和 Stearns[1966]。关于非确定层次性问题, 习题 11.5(空间)由 Ibarra[1972]给出, 习题 11.6 由 Cook[1973]给出。关于 RAM 的层次性问题, 习题 11.11 由 Cook 和 Reckhow[1973]给出。习题 11.14 取自 J. Hopcroft, 习题 11.19 取自 Meyer 和 Stockmeyer[1973]。

对于一些“自然的”问题的指数复杂度下界的工作开始于 Meyer[1972]以及 Meyer 和 Stockmeyer[1972]。关于半扩展正则表达式的定理 11.2 可见 Hunt[1973b]。关于扩展正则表达式的定理 11.3 可见 Meyer 和 Stockmeyer[1973]。其他关于难解问题的可见 Book[1972], Hunt[1973a、1974], Stockmeyer 和 Meyer[1973], Meyer[1972], Hunt 和 Rosenkrantz[1974], Rounds[1973]以及 Constable、Hunt 和 Sahni[1974]。

第12章 算术运算的下界

设计求解一个给定问题的算法时最基本的问题是“该问题固有的计算复杂度是多少?”了解算法在有效性方面的理论下界,可以评估所设计的算法的好坏,以及应再付出多大的努力,以找出更好的解决方案。假设已经知道该问题是难解的,则应该采用启发式方法来寻找合适的解决方案。

遗憾的是,确定某个问题的固有计算复杂度经常是非常困难的。对于大多数实际问题,往往依赖经验确定算法的好坏。但是在某些问题中,可以通过执行某些计算需要的算术操作的次数来确定算法的紧致下界。本章将分析该本质问题的一些基本结论。例如, $n \times p$ 矩阵乘以 p 维向量需要 np 次标量乘法操作,而求解 n 次多项式的结果需要 n 次乘法操作。一些关于下界运算的其他结论包含在习题中。对下界感兴趣的读者应把习题作为本章不可缺少的一部分。

12.1 域

为了获得精确的下界,必须知道允许的基本操作。为了计算的明确性,假定所有的运算都在某个域中完成,比如实数域,其中的基本操作为域中元素的加法和乘法运算。

定义 代数系统 $(A, +, \cdot, 0, 1)$ 为域,必须满足

1. 该系统为环,具有乘法单位元 1。
2. 乘法运算符合交换律。
3. $A - \{0\}$ 中的任意元素 a 存在乘法逆元 a^{-1} ,使得 $aa^{-1} = a^{-1}a = 1$ 。[⊖]

例 12.1 采用普通加法及乘法运算的实数构成域。整数形成环,但不是域,因为除了 ± 1 外,其他元素不存在乘法逆元。但当 p 是素数时,整数模 p 的运算则构成(有限)域。□

给定 x 的值,考虑计算任意多项式 $p(x) = \sum_{i=0}^n a_i x^i$ 结果的问题。希望该算法以所有 a_i 和 x 的值为输入,计算出多项式 $p(x)$ 的相应结果。该算法针对某域的所有可能的输入值进行。算法对所有允许的输入执行的加法、减法及乘法的最大操作次数称为该算法的算术复杂度。

注意,某些 n 次多项式的运算比其他多项式计算要容易。例如,计算 $x^n + 2$ 仅需要 $\log n$ 次操作,而直观上“随机” n 次多项式需要线性次运算。因此,只用于某些特定多项式计算的算法比通用多项式计算算法快。为了得到能够解决各类多项式的算法,在这里用未定元表示所有的输入变量。

定义 代数系统的未定元是符号,而不是确定集合中的元素。令 $F = (A, +, \cdot, 0, 1)$ 为域,且 x_1, x_2, \dots, x_n 为 F 的未定元。未定元 x_1, x_2, \dots, x_n 对 F 的扩展记为 $F[x_1, \dots, x_n]$,它是满足 B 包含 $A \cup \{x_1, \dots, x_n\}$ 的最小交换[⊕]环 $(B, +, \cdot, 0, 1)$ 。

注意,“隐藏”单元与未定元之间是没有关系的,因此,每个多项式的系数在 F 中,且“未知”的元素 x_1, x_2, \dots, x_n 表示 $F[x_1, \dots, x_n]$ 中的某些元素。如果运用交换环的运算定律,两个多项式能够互相转换,则两个多项式表示 $F[x_1, \dots, x_n]$ 的同一个元素。 F 的乘法单元 1,也是 $F[x_1, \dots, x_n]$ 上的乘法单元。

例 12.2 令 F 为实数域,则环 $F[x, y]$ 包括 x, y 及所有的实数。因为 $F[x, y]$ 在 $+$ 上是封闭的,如 $x+y$ 及 $x+4$ 都在 $F[x, y]$ 中。由于 $F[x, y]$ 在乘法上是封闭的,则它包括 $(x+y)(x+$

⊖ 通常情况下,在不引起混淆时略掉。

⊕ 即乘法符合交换律。

4), 根据环的分配律, 它等价于 $x^2 + xy + 4x + 4y$. □

12.2 再论直线状代码

再次考虑计算任意多项式需要多少次算术运算的问题。我们真正关心的是构造表达式 $\sum_{i=0}^n a_i x^i$ 或由未知元 a_0, \dots, a_n 和 x 构造等价的表达式需要多少次 $+$ 及 \cdot 操作。基于这些观察可以建立以下本质上是 1.5 节类似的直线状程序模型的计算模型。

定义 令 F 为一个域, F 上的计算涉及以下几个方面。

1. 称为输入的未定元集合。

2. 变量名集合。

3. 形如 $a \leftarrow b \theta c$ 的操作步骤序列, 其中 θ 为 $+$ 、 $-$ 或 $*$ 等操作, a 是在前面步骤中没有出现过的变量, b 和 c 为输入或为 F 中的元素, 或为在前面步骤中出现在箭头左边的变量名。

为了方便起见, 把 $a \leftarrow b + 0$ 写成 $a \leftarrow b$, 把 $a \leftarrow 0 - b$ 写成 $a \leftarrow -b$ 。计算中出现的 F 中的元素称为常量。

为了确定计算结果, 必须知道计算过程中变量的值。计算是一步一步地完成的, 且每一步把 $F[x_1, \dots, x_n]$ 中的一个元素赋给一个新的变量, 其中 x_i 是输入。对变量或输入 a 的值 $v(a)$ 可定义如下。如果 a 是输入或者 F 中的元素, 则 $v(a) = a$ 。如果 a 是变量且 $a \leftarrow b \theta c$ 是 a 在左边的一步, 则 $v(a) = v(b) \theta v(c)$ 。

计算 E , 基于域 F 存在 $F[x_1, \dots, x_n]$ 的表达式集合, 对于 E 中的任意元素 e , 存在变量 f , 使得 $v(f) = e$ 。

可以看出, 计算是基于相关的域进行的。例如, 计算 $x^2 + y^2$, 此时 F 是实数域, 尽管乘以常数时乘法运算次数不计数, 也需要两次乘法运算。如果 F 是复数域, 则只需要一次乘法运算, 即 $(x + iy)(x - iy)$ (不包括乘以常数)。通常默认域为实数域, 尽管也可以用复数域、有理数域或其他有限域。使用的基本操作决定所选取的域。假设可以表示实数, 对实数执行加法和乘法运算作为基本操作时, 则 F 为实数域。

例 12.3 回忆计算两个复数 $a + ib$ 和 $c + id$ 的乘法运算, 在实数域下, 可表示为两个表达式的计算: $ac - bd$ 和 $ad + bc$ 。计算过程如下:

$$\begin{aligned} f_1 &\leftarrow a * c \\ f_2 &\leftarrow b * d \\ f_3 &\leftarrow f_1 - f_2 \\ f_4 &\leftarrow a * d \\ f_5 &\leftarrow b * c \\ f_6 &\leftarrow f_4 + f_5 \end{aligned}$$

变量 f_1 的值为 ac 。用类似的方法可得到 $v(f_2) = bd$ 和 $v(f_3) = ac - bd$ 。因此, f_3 的值等于第一个表达式。 f_6 的值为 $ad + bc$, 是第二个表达式。这个计算过程可得到两个复数的积。

还有另一种计算两个复数乘法的方法, 仅需要三次实数乘法操作:

$$\begin{aligned} f_1 &\leftarrow a + b \\ f_2 &\leftarrow f_1 * c \\ f_3 &\leftarrow d - c \\ f_4 &\leftarrow a * f_3 \\ f_5 &\leftarrow f_4 + f_2 \\ f_6 &\leftarrow d + c \end{aligned}$$

$$f_7 \leftarrow b * f_6$$

$$f_8 \leftarrow f_2 - f_7$$

容易看出 $v(f_5) = ad + bc$ 且 $v(f_8) = ac - bd$ 。□

明显看出, 当且仅当针对表达式的输入用域 F 中的任意值代替后得到的是该表达式的特定实例, 然后不断利用输入值取代中间未定元得到结果。

本章的后面部分主要考虑计算一系列表达式的乘法操作次数的下界。考虑乘法的原因是在某些领域乘法操作的代价比加法及减法大得多。

例如, 在第6章中看到, 可以用7次标量乘法操作完成两个 2×2 矩阵的乘法, 因此得到 $O_A(n^{\log 7})$ 矩阵乘法算法。事实上, Strassen 算法用18次加法操作在渐近复杂度中可忽略不计。在 Strassen 算法的递归的第一层中, 7次 $(n/2) \times (n/2)$ 矩阵乘法操作所需的时间远比18次(或其他次)同类型矩阵的加减法操作开销大(n 趋于无穷)。

事实上, 如果采用不可交换环的运算规律可在6次标量乘法操作内完成两个 2×2 矩阵乘法, 那么不管 2×2 矩阵乘法在计算过程中用了多少次加法及减法操作, 都可以得到 $O_A(n^{\log 6}) = O_A(n^{2.59})$ 的矩阵乘法算法。(已经证明, 假设矩阵乘法算法对任意环适用, 则 2×2 矩阵乘法需要7次乘法操作。见 Hopcroft 及 Kerr[1971]。)

另一个实例为2.6节中的两个 n 位整数的乘法可以在 $O_B(n^{1.59})$ 时间内完成, 因为3次乘法操作就能完成表达式 ac , bd 和 $ad + bc$ 的计算。事实上, 如果只用2次乘法操作就能计算这些表达式的值, 则存在常数 c , 使得 $M(n) \leq 2M(n/2) + cn$, 递归应用这样的计算过程, 可得到 $O_B(n \log n)$ 的整数乘法算法, 该算法比已知的更好。遗憾的是, 在一些域中, 这些表达式的乘法操作不能低于3次。

12.3 问题的矩阵表述

许多问题可以通过计算矩阵与列向量乘积进行形式化。矩阵的元素来自 $F[a_1, \dots, a_n]$, 其中 F 为域, a_1, \dots, a_n 为中间未知元, 列向量中的各个元素的中间未知元与 a_1, \dots, a_n 不同。

例 12.4 两个复数 $a + ib$ 和 $c + id$ 的乘法可以采用矩阵 - 向量积的方式表示

$$\begin{pmatrix} a & -b \\ b & a \end{pmatrix} \begin{pmatrix} c \\ d \end{pmatrix} = \begin{pmatrix} ac - bd \\ bc + ad \end{pmatrix}$$

这里 F 是实数域, 矩阵里的元素选自 $F[a, b]$ 。向量的未知元为 c 和 d 。□

例 12.5 多项式 $\sum_{i=0}^n a_i x^i$ 的计算可以表示为

$$(1, x, x^2, \dots, x^n) \begin{pmatrix} a_0 \\ a_1 \\ \vdots \\ a_n \end{pmatrix}$$

这里 F 是实数域, $1 \times (n+1)$ 矩阵的元素来自 $F[x]$ 。□

12.4 面向行的矩阵乘法的下界

定义 令 F 是域, 且 a_1, \dots, a_n 为中间未知元, 用 $F^m[a_1, \dots, a_n]$ 表示 m 维的向量空间, 其中的元素来自 $F[a_1, \dots, a_n]$ 。令 F^m 为 m 维的向量空间, 其中元素来自 F 。来自 $F^m[a_1, \dots, a_n]$ 的向量集合 $\{v_1, \dots, v_k\}$ 是模 F^m 线性无关 (linearly independent modulo F^m) 的, 如果对 F 中的 u_1, \dots, u_k , $\sum_{i=1}^k u_i v_i$ 在 F^m 中, 蕴涵所有的 u_i 为零。如果 v_i 不是模 F^m 线性无关的, 称该集合模 F^m 相关 (dependent modulo F^m)。

另一种检查模 F^m 线性无关的方法是考虑在 $F^m[a_1, \dots, a_n]$ 中的向量的等价类的商空间 $F^m[a_1, \dots,$

$a_n]/F^m$ (当且仅当 $v_1 - v_2$ 在 F^m 中时, 向量 v_1, v_2 在 $F^m[a_1, \dots, a_n]$ 中等价)。模 F^m 线性无关意味着 $F^m[a_1, \dots, a_n]/F^m$ 的等价类线性无关。

例 12.6 假定 F 为实数域, 且 a 和 b 为中间未知元, 则以下三个向量

$$v_1 = \begin{pmatrix} 2a \\ 4b \\ 2a - 2b \end{pmatrix}, \quad v_2 = \begin{pmatrix} -a \\ -b \\ b + 3 \end{pmatrix}, \quad v_3 = \begin{pmatrix} \pi \\ b - 1 \\ a \end{pmatrix}$$

在 $F^3[a, b]$ 中且模 F^3 相关, 因为

$$\frac{1}{2}v_1 + v_2 - v_3 = \begin{pmatrix} -\pi \\ 1 \\ 3 \end{pmatrix}$$

另一方面, 在 F^3 中, 向量

$$v_1 = \begin{pmatrix} b \\ a \\ 0 \end{pmatrix}, \quad v_2 = \begin{pmatrix} a \\ 0 \\ b \end{pmatrix}, \quad v_3 = \begin{pmatrix} 0 \\ b \\ a \end{pmatrix}$$

是模 F^2 线性无关。假定 $\sum_{i=1}^3 u_i v_i$ 在 F^2 中, 则 $u_1 b + u_2 a$ 在 F 中。由于 a 和 b 都不是 F 中的元素, 则必须有 $u_1 = u_2 = 0$ 。另外, 由 $u_1 a + u_3 b$ 在 F 中, 可以得到 $u_1 = u_3 = 0$ 。最后得出 v_1, v_2, v_3 是模 F^2 线性无关。□

如果 M 是 $r \times p$ 矩阵, 且矩阵元素来自 F , 则 M 中线性无关的行的数目等于 M 中线性无关的列的数目。但是, 如果 M 中的元素来自 $F[a_1, \dots, a_n]$, 则模 F^p 线性无关的行的数目不一定等于模 F^r 线性无关的列的数目。例如, 1×3 矩阵 $[a_1, a_2, a_3]$ 具有一个模 F^3 线性无关的行和 3 个模 F 线性无关的列。矩阵 M 模 F^p 的行秩是模 F^p 线性无关的行的数目。矩阵 M 模 F^r 的列秩是模 F^r 线性无关的列的数目。

以下的术语将在本节后面部分及后面两节经常使用。

1. F 表示域。
2. $\{a_1, \dots, a_n\}$ 和 $\{x_1, \dots, x_p\}$ 表示不相交的中间未知元的集合。
3. M 是 $r \times p$ 矩阵, 且元素来自 $F[a_1, \dots, a_n]$ 。
4. x 是列向量 $[x_1, \dots, x_p]^T$ 。[⊖]

定理 12.1 令 M 为元素来自 $F[a_1, \dots, a_n]$ 的矩阵, 且令 x 为列向量 $[x_1, \dots, x_p]^T$ 。假定 M 的行秩为 r , 则计算 Mx 至少需要 r 次乘法操作。

证明: 不失一般性, 假定 M 有 r 行 (否则, 通过删掉多余的行, 使之只剩下 r 个线性无关的行。令 M' 为最终的矩阵, 所有 Mx 的计算都用 $M'x$ 表示。最后可得到, 如果 $M'x$ 需要 r 次乘法操作, 则 Mx 也需要 r 次乘法操作)。

假定在计算 Mx 的过程中需要 s 次乘法操作。令 e_1, \dots, e_s 为每次乘法操作的表达式, 则 Mx 的元素可以表示为 e_1, \dots, e_s 的线性函数与未知元的和。可以写出

$$Mx = Ne + f \quad (12-1)$$

其中 e 为向量, 且向量的第 i 个元素为 e_i , f 为向量, 其元素为 a_i 和 x_i 的线性函数。 N 为 $r \times s$ 的矩阵, 元素来自 F 。

现在假定 $r > s$ 。考虑线性代数里的基本结论: 如果 $r > s$, 则 N 的行线性相关[⊖]。也就是说在 F^r 中存在 $y \neq 0$ [⊖], 使得 $y^T N = 0^T$ 。式 (12-1) 乘以 y^T , 可得:

⊖ 以列的形式给出列向量是不方便的, 因此和第 7 章一样, 以行的方式表示列向量, 并用上标 T 表示转置。

⊖ 如果读者不熟悉该结论, 即线性代数的基本定理, 引理 12.1 叙述的性质则有着更强的结论, 其证明过程足以作为证明当前结论所需的线索。

⊖ 0 是维数适当的分量全 0 的向量。

$$\mathbf{y}^T M \mathbf{x} = \mathbf{y}^T N \mathbf{e} + \mathbf{y}^T \mathbf{f} = \mathbf{y}^T \mathbf{f} \quad (12-2)$$

从式(12-2)可以看出, $(\mathbf{y}^T M) \mathbf{x} = \mathbf{y}^T \mathbf{f}$, 它是中间未知元的线性函数。由于 \mathbf{x} 在每个分量中具有不同的未定元, 因此, 可以推断出行向量 $\mathbf{y}^T M$ 中的元素仅来自 F 。否则, 将有一个项涉及未定元在 $\mathbf{y}^T M \mathbf{x}$ 中的乘法结果, 因此为 $\mathbf{y}^T \mathbf{f}$, 这与 \mathbf{f} 是未定元的线性假定相矛盾。由于 $\mathbf{y} \neq \mathbf{0}$ 且 $\mathbf{y}^T M$ 在 F^p 中, M 的依赖模 F^p 相关, 也与假设相矛盾。因此 $s \geq r$ 且定理成立。□

例 12.7 12.2 节曾考虑为一个简单的整数递归乘法算法计算三个表达式 ac , bd 和 $ad + bc$ 。这三个公式可以表示为矩阵 - 向量的乘法:

$$\begin{pmatrix} a & 0 \\ 0 & b \\ b & a \end{pmatrix} \begin{pmatrix} c \\ d \end{pmatrix} \quad (12-3)$$

令 \mathbf{R} 为实数域。例 12.6 说明式(12-3)中矩阵的三行是模 \mathbf{R}^2 线性无关的。因此, 计算以上关于实数的三个表达式需要三次实数乘法操作。□

12.5 面向列的矩阵乘法的下界

类似定理 12.1 中的线性无关列的结果在这里也是成立的, 但比较难获得。在此, 需要一些线性无关向量的基本结论。

引理 12.1 令 $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$ 为 m 维向量的集合, 其中 $k \geq 1$, 且元素来自 $F[a_1, \dots, a_n]$ 。假定 $\{\mathbf{v}_i \mid 1 \leq i \leq k\}$ 有 q 个向量的子集, 且这些向量是模 F^m 线性无关的, 则对于 F 中的任意元素 b_2, \dots, b_k , 集合 $\{\mathbf{v}'_i \mid \mathbf{v}'_i = \mathbf{v}_i + b_i \mathbf{v}_1, 2 \leq i \leq k\}$ 包含 $q-1$ 个向量的子集, 且这些向量是模 F^m 线性无关的。

证明: 对 $\mathbf{v}_2, \mathbf{v}_3, \dots, \mathbf{v}_k$ 进行重新编号, 如果需要, 可以假设 $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_q\}$ 或者 $\{\mathbf{v}_2, \mathbf{v}_3, \dots, \mathbf{v}_{q+1}\}$ 线性无关。○

情况 1 假定 $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_q\}$ 是线性无关的, 则得出 $\{\mathbf{v}'_2, \mathbf{v}'_3, \dots, \mathbf{v}'_q\}$ 是线性无关的。为说明这一点, 假定对于 F 中 c_i 的某个集合, $\sum_{i=2}^q c_i \mathbf{v}'_i$ 在 F^m 中, 则当 $c_1 = \sum_{i=2}^q c_i b_i$ 时, $\sum_{i=1}^q c_i \mathbf{v}_i$ 在 F^m 中。但是, 假如 $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_q\}$ 是线性无关的, 则 c_i 为 0, $1 \leq i \leq q$ 。因此 $\{\mathbf{v}'_2, \mathbf{v}'_3, \dots, \mathbf{v}'_q\}$ 是线性无关的。

情况 2 假定 $\{\mathbf{v}_2, \mathbf{v}_3, \dots, \mathbf{v}_{q+1}\}$ 是线性无关的。如果 $\{\mathbf{v}'_2, \mathbf{v}'_3, \dots, \mathbf{v}'_q\}$ 是线性无关的, 则可得结果; 如果不是线性无关的, 则 F 中存在不全为 0 的 c_2, \dots, c_q , 使得 $\sum_{i=2}^q c_i \mathbf{v}'_i$ 在 F^m 中。令 $c_1 = \sum_{i=2}^q c_i b_i$, 则 $\mathbf{w} = \sum_{i=1}^q c_i \mathbf{v}_i$ 也在 F^m 中。

可以假定 $c_1 \neq 0$, 否则 $\sum_{i=2}^q c_i \mathbf{v}_i$ 将在 F^m 中, 与 $\{\mathbf{v}_2, \mathbf{v}_3, \dots, \mathbf{v}_{q+1}\}$ 线性无关的假定相矛盾。因此, 可以写出下式:

$$\mathbf{v}_1 = c_1^{-1} \left(\mathbf{w} - \sum_{i=2}^q c_i \mathbf{v}_i \right)$$

由于 c_2, \dots, c_q 不全都为 0, 因此, 不失一般性, 我们假设 $c_2 \neq 0$ 。我们再次使用上述的论证方法对集合 $\{\mathbf{v}'_3, \mathbf{v}'_4, \dots, \mathbf{v}'_{q+1}\}$ 进行证明。假设该集合是线性无关的, 则可以假设 $\sum_{i=3}^{q+1} d_i \mathbf{v}'_i$ 在 F^m 中, 其中 d_i 在 F 中且不全为 0。令 $d_1 = \sum_{i=3}^{q+1} d_i b_i$, 则 $\mathbf{x} = d_1 \mathbf{v}_1 + \sum_{i=3}^{q+1} d_i \mathbf{v}_i$ 在 F^m 中。如果 $d_1 = 0$, 则与 $\{\mathbf{v}_3, \dots, \mathbf{v}_{q+1}\}$ 线性无关的假设相矛盾。如果 $d_1 \neq 0$, 则可以写出下式:

$$\mathbf{v}_1 = d_1^{-1} \left(\mathbf{x} - \sum_{i=3}^{q+1} d_i \mathbf{v}_i \right)$$

○ 整个证明略去了“模 F^m ”。

由这两个关于 v_i 的式子相等, 可得:

$$d_1^{-1} \left(\mathbf{x} - \sum_{i=3}^{q+1} d_i v_i \right) = c_1^{-1} \left(\mathbf{w} - \sum_{i=2}^q c_i v_i \right)$$

两边同时乘以 $c_1 d_1$ 并对各项进行重新排列可得

$$d_1 c_2 v_2 + \sum_{i=3}^q (d_1 c_i - c_1 d_i) v_i - c_1 d_{q+1} v_{q+1} = d_1 \mathbf{w} - c_1 \mathbf{x}$$

由于 \mathbf{w} 和 \mathbf{x} 在 F^m 中, 因此 $d_1 \mathbf{w} - c_1 \mathbf{x}$ 也在 F^m 中。由于 c_2 和 d_1 为非零数, 则 $\{v_2, v_3, \dots, v_{q+1}\}$ 线性相关, 与假设矛盾。□

令 M 为矩阵元素来自 $F[a_1, \dots, a_n]$ 的 $r \times p$ 矩阵, 现在证明, 如果 M 的 q 列是模 F^r 线性无关的 ($q \geq 1$), 则 $M\mathbf{x}$ 的任意的计算过程至少需要 q 次乘法操作, 其中 $\mathbf{x} = [x_1, \dots, x_p]^T$ 。

定义: 如果一个乘数值涉及未定元 x_i , 另一个乘数不是 F 中的元素, 则该乘法称为活跃的。例如, 如果 $v(b) = 3 + a_2$, $v(c) = x_1 + 2 * x_3$, 则称 $b * c$ 为活跃的。但是, 如果 $v(b) = 3$ 且 $v(c) = x_1 + 2 * x_3$, 或者 $v(b) = 3 + a_2$ 且 $v(c) = a_1 + 2 * a_3$, 则 $b * c$ 就不是活跃的。

定理 12.2 令 y 为元素来自 $F[a_1, \dots, a_n]$ 的 r 维向量。如果 M 的列秩为 q , 则计算过程 $M\mathbf{x} + y$ 至少需要 q 次活跃乘法, 其中 $q \geq 1$ 。

证明: 通过对 q 进行归纳来证明。

归纳基础: $q=1$ 。在 M 中存在一列不在 F^r 中, 因此 M 中存在元素 e , 它在 $F[a_1, \dots, a_n]$ 中, 但不在 F 中。因此, $M\mathbf{x}$ 中的某个元素对某个 j , 具有项 ex_j , $M\mathbf{x} + y$ 也同样成立。没有活跃乘法的计算过程只能计算形如 $P(a_1, \dots, a_n) + L(x_1, \dots, x_p)$ 的表达式, 其中 P 是多项式, L 是一个线性函数, 每个系数都在 F 中。因此, 计算 $M\mathbf{x} + y$ 至少有一次活跃乘法。

归纳步骤: 假定 $q > 1$ 且定理对 $q-1$ 成立。令 C 为 $M\mathbf{x} + y$ 的计算过程。由归纳假设可知, C 至少有 $q-1 \geq 1$ 次活跃乘法。假定 $f \leftarrow g * h$ 是第一个这样的乘法, 则不失一般性, 假定

$$v(g) = P(a_1, \dots, a_n) + \sum_{i=1}^p c_i x_i \quad (12-4)$$

其中 P 是系数在 F 中的多项式, 且 c_i 在 F 中。更进一步, 不失一般性, 假设 $c_1 \neq 0$ 。

这时, 策略是通过从 C 和表达式 $M\mathbf{x} + y$ 的集合中构造新的表达式集合 $M'\mathbf{x}' + y'$ 形成计算过程 C' , C' 具有比 C 更少的活跃乘法。进一步, M' 的 $q-1$ 列将是模 F^r 线性无关的。因此, 通过归纳假设, C' 具有 $q-1$ 次活跃乘法, 意味着 C 有 q 次活跃乘法。特别地, 把 C 中构成 g 的表达式 e 中的 x_i 替换成 0, 而且表达式 e 可在没有活跃乘法时计算出来。计算过程 C' 从计算 e 开始, 表达式 e 的值赋予 x_1 (x_1 将不再作为输入变量)。 C' 的剩余部分由原来的 C 中的 $f \leftarrow g * h$ 替换为 $f \leftarrow 0$, 其他保持不变。因此, C' 计算的表达式可以表示为 $M'\mathbf{x}' + y'$, 这里 M' 的 $q-1$ 列是模 F^r 线性无关的。

现在回到证明细节中, 从式(12-4)和假设 $c_1 \neq 0$, 可得 $g=0$, 条件是下式成立

$$x_1 = -c_1^{-1} \left[P(a_1, \dots, a_n) + \sum_{i=2}^p c_i x_i \right] \quad (12-5)$$

式(12-5)的右边是上面提到的表达式 e 。通过前面描述的 C 形成的计算过程 C' 通过把 e 替换成 x_1 来计算 $M\mathbf{x} + y$, 因此 $M'\mathbf{x}' + y'$ 可以写成

$$M \begin{bmatrix} e \\ x_2 \\ x_3 \\ \vdots \\ x_p \end{bmatrix} + y$$

可以进一步写成

$$M \begin{bmatrix} -c_1^{-1} \sum_{i=2}^p c_i x_i \\ x_2 \\ x_3 \\ \vdots \\ x_p \end{bmatrix} + M \begin{bmatrix} -c_1^{-1} P(a_1, \dots, a_n) \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} + y \quad (12-6)$$

考虑式(12-6)中的第一项, 令 M 的第 i 列为 \mathbf{m}_i 。对于 $2 \leq i \leq p$, 定义 $\mathbf{m}'_i = \mathbf{m}_i - (c_i/c_1) \mathbf{m}_1$ 。令 M' 为 $p-1$ 列矩阵, 其第 i 列为 \mathbf{m}'_{i+1} 且令 $\mathbf{x}' = [x_2, \dots, x_p]^T$ 。因此, 式(12-6)中的第一项等于 $M' \mathbf{x}'$ 。

式(12-6)的第二项是元素来自 $F[a_1, \dots, a_n]$ 的向量, 因为矩阵 M 中的元素在 $F[a_1, \dots, a_n]$ 中。因此式(12-6)中的第2和第3项能够组合形成新的向量 \mathbf{y}' , 其元素来自 $F[a_1, \dots, a_n]$ 。因此, C' 计算 $M' \mathbf{x}' + \mathbf{y}'$ 。从引理 12.1 可以直接推出 M' 至少有 $q-1$ 列是模 F' 线性无关的。

因此, C' 具有 $q-1$ 次活跃乘法, 表明 C 具有 q 次活跃乘法。□

这里给出运用定理 12.2 的两个例子。

例 12.8 把 $n \times p$ 的矩阵 A 乘以长度为 p 的向量 \mathbf{v} 需要 np 次乘法。形式地说, 对于 $1 \leq i \leq n$ 和 $1 \leq j \leq p$, 令 a_{ij} 和 v_j 为未定元, 则 $A\mathbf{v}$ 是如图 12-1 所示的矩阵-向量乘法。

$$\begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1p} \\ a_{21} & a_{22} & \cdots & a_{2p} \\ \vdots & \vdots & & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{np} \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_p \end{bmatrix}$$

图 12-1 一般的矩阵
向量乘法

直接对 $A\mathbf{v}$ 应用定理 12.1 或 12.2, 可得计算只需要 $\text{MAX}(n, p)$ 次乘法操作。然而, 矩阵-向量乘法操作同样也可以表示为 $M\mathbf{x}$, 其中 M 的第 i 行具有在 v_1, v_2, \dots, v_p 中的 $(i-1)p+1$ 列 $\sim ip$ 列有值, 其他列为 0。向量 \mathbf{x} 是列向量, 由 A 的各行连接起来形成。 M 和 \mathbf{x} 如图 12-2 所示。

$$\begin{array}{c} n \text{ 行} \\ \begin{bmatrix} v_1 & v_2 & \cdots & v_p & 0 & 0 & \cdots & 0 & \cdots & 0 & 0 & \cdots & 0 \\ 0 & 0 & \cdots & 0 & v_1 & v_2 & \cdots & v_p & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & \vdots & & \vdots & \vdots & \vdots & & \vdots & & \vdots & \vdots & & \vdots \\ 0 & 0 & \cdots & 0 & 0 & 0 & \cdots & 0 & \cdots & v_1 & v_2 & \cdots & v_p \end{bmatrix} \end{array} \quad \begin{array}{c} \begin{bmatrix} a_{11} \\ a_{12} \\ \vdots \\ a_{1p} \\ a_{21} \\ a_{22} \\ \vdots \\ a_{2p} \\ \vdots \\ a_{n1} \\ a_{n2} \\ \vdots \\ a_{np} \end{bmatrix} \end{array}$$

np 列

图 12-2 矩阵-向量乘法等价形式

很容易看出, M 的列是模 F^n 线性无关的。通过定理 12.2, $A\mathbf{v}$ 的计算过程至少需要 np 次乘法操作。□

例 12.9 考虑 n 次多项式 $\sum_{i=0}^n a_i x^i$ 的计算问题, 该问题可以表示为矩阵-向量乘法 $M\mathbf{x}$:

$$[1, x, x^2, \dots, x^n] \begin{bmatrix} a_0 \\ a_1 \\ \vdots \\ a_n \end{bmatrix}$$

这里 M 中元素来自 $F[x]$ 且 $\mathbf{X} = [a_0, a_1, \dots, a_n]^T$ 。可以很容易地看出, M 中列的集合(除了第一列以外)形成一个模线性无关的集合。因此, M 具有 n 个线性无关的列, 计算任意的 n 次多项式需要 n 次乘法操作。

霍纳规则采用以下方案计算多项式

$$(\dots((a_n x + a_{n-1})x + a_{n-2})x + \dots + a_1)x + a_0$$

它需要 n 次乘法操作, 且是经过优化的, 需要尽可能少的乘法。采用类似的方式可以得到计算 $\sum_{i=0}^n a_i x^i$ 时需要 n 次加法或者减法。因此, 在某点计算多项式时运用霍纳规则所需的算术运算次数是最少的。□

12.6 面向行和列的矩阵乘法下界

把定理 12.1 和 12.2 结合起来可获得比单独考虑行和列时更强的结果。令 M 为 $r \times p$ 矩阵, 元素来自 $F[a_1, \dots, a_n]$, $\mathbf{x} = [x_1, \dots, x_p]^T$ 。

定理 12.3 假定 M 具有子矩阵 S , S 有 q 行和 c 列。对 F^q , F^c 中的任意向量 \mathbf{u} 和 \mathbf{v} , 当且仅当 $\mathbf{u} = \mathbf{0}$ 或 $\mathbf{v} = \mathbf{0}$ 时, $\mathbf{u}^T S \mathbf{v}$ 是 F 中的元素, 则任意计算 $M\mathbf{x}$ 的过程至少需要 $q+c-1$ 次乘法操作。

证明: 不失一般性, 假设 M 本身只有 q 行, 且 S 由 M 的前 c 列组成。假定 $M\mathbf{x}$ 可以在 s 次乘法操作内完成。令 \mathbf{e} 为元素来自从乘法得到的 s 表达式的向量。假定 \mathbf{e} 中第 i 个元素在第 j 个元素前计算, 且 $i < j$, 则定理 12.1 可以写成

$$M\mathbf{x} = N\mathbf{e} + \mathbf{f} \quad (12-7)$$

其中 N 是元素来自 F 的 $q \times s$ 矩阵, \mathbf{f} 是向量, 它的元素是 a_i 和 x_i 的线性组合。

根据定理 12.1, 可以得出 $s \geq q$ 。如果不是, 可以找出 F^q 中的一个向量 $\mathbf{y} \neq \mathbf{0}$, 使 $\mathbf{y}^T N = \mathbf{0}^T$, 表明 $\mathbf{y}^T M$ 中的元素在 F 中。因此, $\mathbf{y}^T S$ 中的元素也在 F 中, 与假设矛盾(取 $\mathbf{u} = \mathbf{y}$ 及 $\mathbf{v} = [1, \dots, 1]$)。

由于 $s \geq q$, 可以把 N 分成两个矩阵 N_1 和 N_2 , 其中 N_1 由 N 的前 $s-q+1$ 列构成, N_2 由 N 的后 $q-1$ 列构成。同样, 令 \mathbf{e}_1 和 \mathbf{e}_2 分别为 \mathbf{e} 的前 $s-q+1$ 个元素和 $q-1$ 个元素。可以把式(12-7)写成

$$M\mathbf{x} = N_1 \mathbf{e}_1 + N_2 \mathbf{e}_2 + \mathbf{f} \quad (12-8)$$

由于 N_2 为 $q \times (q-1)$, 则存在 F^q 中的向量 $\mathbf{y} \neq \mathbf{0}$, 使 $\mathbf{y}^T N_2 = \mathbf{0}^T$ 。令式(12-8)乘以 \mathbf{y}^T , 产生

$$\mathbf{y}^T M\mathbf{x} = \mathbf{y}^T N_1 \mathbf{e}_1 + \mathbf{y}^T \mathbf{f} \quad (12-9)$$

令 $M' = \mathbf{y}^T M$ 。注意, M' 是个 $1 \times p$ 矩阵, 它是 M 各行的线性组合。由于 \mathbf{e}_1 中的结果可以在不涉及 \mathbf{e}_2 的情况下计算(前者假定先计算), 显然可以得出 $M'\mathbf{x}$ 可采用式(12-9)并花费 $s-q+1$ 次乘法操作完成。如果可以证明 M' 有 c 列是模 F 线性无关的, 则由定理 12.2, 有 $s-q+1 \geq c$, 即 $s \geq q+c-1$, 得证。

现在证明 $M' = \mathbf{y}^T M$ 的前 c 列是模 F 线性无关的。令 $\mathbf{y}^T = [y_1, \dots, y_q]$, M' 的前 c 项为 $\sum_{i=1}^q y_i M_{ij}$ ($1 \leq j \leq c$), M_{ij} 为 M 的第 j 个元素。假定存在向量 $\mathbf{z} \neq \mathbf{0}$, 它的元素为 z_1, \dots, z_c , 且 z_i 在 F 中, 使得 $\sum_{j=1}^c z_j \sum_{i=1}^q y_i M_{ij}$ 在 F 中。也就是说, M' 的前 c 列是模 F 相关的, 则 $\mathbf{y}^T S \mathbf{z}$ 在 F 中, 与对 S 的假设相矛盾。因此, 可得 M' 有 c 列是模 F 线性无关的, 定理得证。□

现在将定理 12.3 应用到复数 $a+ib$ 和 $c+id$ 的乘法, 以表明该运算需要 3 次实数乘法操作。由观察可得, 单独采用定理 12.1 和 12.2 对于证明这个问题来说是不够的。

例 12.10 考虑两个复数 $a+ib$ 和 $c+id$ 相乘的问题, 即计算

$$M\mathbf{x} = \begin{pmatrix} a & -b \\ b & a \end{pmatrix} \begin{pmatrix} c \\ d \end{pmatrix}$$

令 S 为 M 本身。令 F 为实数, 且令

$$\begin{pmatrix} y_1 \\ y_2 \end{pmatrix} \neq \mathbf{0} \quad \text{和} \quad \begin{pmatrix} z_1 \\ z_2 \end{pmatrix} \neq \mathbf{0}$$

在 F^2 中。假定

$$[y_1 \ y_2] \begin{pmatrix} a & -b \\ b & a \end{pmatrix} \begin{pmatrix} z_1 \\ z_2 \end{pmatrix}$$

是 F 中元素, 上述的结果是 $y_1 z_1 a + y_2 z_1 b + y_2 z_2 a - y_1 z_2 b$ 。如果这个表达式是 F 中的元素, 则 a 和 b 的系数必须为 0, 因此

$$y_1 z_1 + y_2 z_2 = 0 \quad (12-10)$$

和

$$y_2 z_1 - y_1 z_2 = 0 \quad (12-11)$$

假定 $y_1 \neq 0$, 则从式(12-11)得到 $z_2 = y_2 z_1 / y_1$, 将它代入式(12-10), 再乘以 y_1 , 得到 $(y_1^2 + y_2^2) z_1 = 0$ 。由于 $y_1 \neq 0$, 可得 $y_1^2 + y_2^2 \neq 0$, 则 $z_1 = 0$ 。通过式(12-11)可得 $z_2 = 0$ 。但是 $z_1 = z_2 = 0$ 与下面的假设矛盾

$$\begin{pmatrix} z_1 \\ z_2 \end{pmatrix} \neq \mathbf{0}$$

现在假定 $y_1 = 0$ 。如果 $y_2 = 0$, 马上可以发现与假设矛盾。如果 $y_2 \neq 0$, 则可调换 y_1 和 y_2 的角色, 可得 $z_1 = z_2 = 0$, 与假设矛盾。

因此, 将定理 12.3 运用到 $M\mathbf{x}$ 中, 其中 $q=c=2$, 需要 3 次实数乘法操作。例 12.3 中的程序表明使用 3 次乘法操作足够了。□

定理 12.1 ~ 12.3 可以推广到在计算中使用除法操作的情景。当在运算步骤中使用除法操作时, 必须允许在步骤中包括某些值除以 0 的操作。因此, 关于乘法操作的理论 (即乘和除) 假定域 F 是无限的。如果 F 是有限的, 可以讨论没有输入的计算过程, 因为每个输入都可能引起除以 0 的操作。

通过上述修改, 定理 12.1 ~ 12.3 可以适用于与乘法有关的操作, 而不只是乘法操作。在定理 12.2 中, 活跃乘法操作的定义 $f \leftarrow g * h$ 或 $f \leftarrow g/h$ 必须 g 和 h 至少有一个值涉及一个且另一个操作数不在 F 中; 或者 g 在 F 中, h 涉及其中的一个 x , 且运算为除法操作。

12.7 预处理

从例 12.9 可以看出, 在一个点计算 n 次多项式需要 n 次乘法。然而, 这里有一个潜在的假设, 就是多项式采用它的系数表示。如果允许使用基于系数计算出来的参数集合来表示多项式, 再来考虑多项式计算所需的最少乘法次数。如果需要多次计算多项式, 发掘不同的多项式表示

⊖ 注意 F 是实数域这个假设在这里非常重要。

方式是提高多项式计算的一种有效方法。这种表示的改变称为预处理。^①在例 12.9 的多项式计算中, 如果乘法运算的一个因子涉及的系数为 a_0, a_1, \dots, a_n , 而另一个因子不在 F 中, 则称该乘法为活跃的。因此, 如果乘法的两个因子数都涉及系数, 但两个因子都不涉及变量 x , 则该乘法可直接计算。通过预处理, 所有没有涉及变量的项都可以直接计算出来。

定义 多变量多项式的次数为多项式中变量最大的次数。例如, $p(x, y) = x^3 + x^2y^2$ 的次数为 3。

下一个引理说明对于任意 $v+1$ 个 v 元多项式, 存在非平凡的关于给定多项式(其值为 0)的多项式函数。例如, 令 $p_1 = x^2, p_2 = x + x^3$, 则对于所有 $x, p_1 + 2p_1^2 + p_1^3 - p_2^2 = 0$ 。

引理 12.2 设 $p_i(x_1, \dots, x_r), 1 \leq i \leq n$, 表示 n 个多项式。如果 $n > r$, 则存在多项式 $g(p_1, \dots, p_n)$, 它的系数不等于 0, 但是作为 x_i 的函数时等于 0。

证明: 令 S_d 为形如 $p_1^{i_1} p_2^{i_2} \dots p_n^{i_n}$ 的项的集合, 其中对于所有的 $j, 0 \leq j \leq d$, 可得 $\|S_d\| = (d+1)^n$ 。令 m 为 p_i 中的最高次数, 则对于每个 $q \in S_d$ 是关于 x_1, x_2, \dots, x_r 的多项式, 且次数至多为 dnm 。

次数为 dnm 的 r 元多项式可以表示为长度是 $(dmn+1)^r$ 的向量。向量中的元素是各项的系数, 且按照固定的顺序排列。因此, 讨论多项式的线性无关性, 特别是基于线性代数的理论, 来自 S_d 中的 $(dmn+1)^r + 1$ 个多项式的集合肯定线性相关即具有一个非平凡的值为 0 的多项式的和。

特别地, 如果 m, n 和 r 为固定的, 且 $n > r$, 对于足够大的 d , 很容易证明 $(d+1)^n \geq (dmn+1)^r$, 因为 $(d+1)^n$ 是一个在 d 上比 $(dmn+1)^r$ 次数更高的多项式。因此, 存在一个 d , 使得 S_d 的成员的非凡和为 x 的函数且等于 0。该和就是所期望的多项式 g 。□

使用引理 12.2, 可得任意参数的集合来表示一个多项式, 则需要 $n+1$ 个参数表示任意的 n 次多项式, 而该多项式的系数是参数的多项式函数。

引理 12.3 设 M 为从 n 维系数向量空间 C^n 到 r 维参数向量空间 D^r 的一对一映射。如果 M^{-1} 使 C^n 中一个向量的每个元素是 D^r 中相应向量的元素的多项式函数, 则 $r \geq n$ 。

证明: 假定 $r < n$ 且给定 $M^{-1}: c_i = p_i(d_1, \dots, d_r)$, 其中 p_i 为多项式, d_i 是参数, c_i 为系数。通过引理 12.2, 存在非平凡多项式 g , 使得对于 d_i 的所有值, $g(c_1, \dots, c_n)$ 等于 0。^② 由于 g 本身不等于 0, 则存在系数为 c_1, \dots, c_n 的多项式, 使得 $g(c_1, \dots, c_n)$ 不是 0。因此, 多项式不能表示为的形式, 矛盾。□

现在说明尽管采用了预处理, 但计算 n 次多项式至少还需要 $n/2$ 次乘法操作。

定理 12.4 即使不计计算过程中仅涉及系数的乘法操作, 任意 n 次多项式在某点的计算过程至少需要 $n/2$ 次乘法操作。

证明: 假定存在一个需要 m 个乘法步骤的涉及未知量 x 的计算过程。令乘法的结果是表达式 f_1, f_2, \dots, f_m , 则每个 f_i 可以表示为

$$f_i = [L_{2i-1}(f_1, \dots, f_{i-1}, x) + \beta_{2i-1}] * [L_{2i}(f_1, \dots, f_{i-1}, x) + \beta_{2i}]$$

其中每个 L_i 都是 f 和 x 的线性函数, 每个 β_i 是系数的多项式函数, 则多项式 $p(x)$ 的计算可以表示为

$$p(x) = L_{2m+1}(f_1, \dots, f_m, x) + \beta_{2m+1}$$

因此, $p(x)$ 可以表示为 $2m+1$ 个参数 β_i 的新集合, 进而, $p(x)$ 的系数是 β_i 的多项式函数。由引理 12.3, 有 $2m+1 \geq n+1$, 因此 $m \geq n/2$ 。□

① 注意这里的情况和 8.5 节的情况不同, 那里关于一个多项式的所有计值点都是已知的。这里仅可以在计值前预处理系数, 而每次进行一次计值, 即计算是在线的而非离线的。

② 注意 g 是系数为 c_i 的多项式, 因为 c_i 是系数为 d_i 的多项式, 因此 g 也是系数为 d_i 的多项式。

定理 12.4 给出一个多项式在某点使用多项式的预处理计算所需要的乘法次数的下界。把多项式重新表示,使其可以在 $n/2$ 次乘法操作内完成的问题具有很大的意义。在这个问题中,不仅要找出容易计算的参数集合,同时还需给出从系数计算出参数的方便方法。因此,我们期望参数是系数的有理函数,而寻找这样的参数的技术会在习题中提及。

习题

12.1 设 F 为域, x 是未知变量。证明变量为 x , 系数来自 F 的多项式集合形成一个交换环 $F[x]$, 而 F 中的乘法单位元是 $F[x]$ 中的乘法单位元。

12.2 证明对以下表达式的计算

$$ac + bd$$

$$ac - bd$$

$$bc + ad$$

$$bc - ad$$

至少需要 4 次乘法操作, 其中 a, b, c 和 d 来自域 F 。

12.3 证明计算列向量与行向量的乘积

$$\begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_m \end{bmatrix} [b_1, b_2, \dots, b_n]$$

至少需要 mn 次乘法操作。

12.4 设二元 x 和 y 的 n 次多项式为 $\sum_{i=0}^n \sum_{j=0}^n a_{ij} x^i y^j$ 。证明对于基于预处理的该多项式的求值, $(n+1)(n+2)/2 - 1$ 次乘法操作是充分且必要的。

12.5 证明 2×2 特普利茨矩阵(见第 6 章的习题)乘以向量

$$\begin{bmatrix} b & c \\ a & b \end{bmatrix} \begin{bmatrix} d \\ e \end{bmatrix}$$

需要 3 次乘法操作。

12.6 对习题 12.5 进行推广证明 $n \times n$ 特普利茨矩阵乘以向量需要 $2n - 1$ 次乘法操作。假设 $n = 3$, 如何得到界限?

*12.7 对 2.6 节中的乘法算法进行推广, 把乘数和被乘数分别划分成 3 块, 然后再用尽量少的乘法, 计算以下矩阵-向量乘积, 一共需要多少次乘法操作? 相应 2.6 节中的算法, 这种方法有改进吗?

$$\begin{bmatrix} a & 0 & 0 \\ b & a & 0 \\ c & b & a \\ 0 & c & b \\ 0 & 0 & c \end{bmatrix} \begin{bmatrix} d \\ e \\ f \end{bmatrix}$$

**12.8 令 F 为域, 令 a_1, \dots, a_n 及 x_1, \dots, x_p 为两个互不相交的未定元集合, 令 $\mathbf{x} = [x_1, \dots, x_p]^T$ 且 M 为 $r \times p$ 矩阵, 其中元素来自 $F[a_1, \dots, a_n]$, 具有 q 个模 F 线性无关的列。证明对 \mathbf{x} 进行预处理, 例如, 不对仅涉及 x_i 的项进行计算, 说明计算 $M\mathbf{x}$ 仍然需要 $q/2$ 次乘法操作。

**12.9 假定集合 P 恰含有 k 个单乘法表达式, 不存在是常数的线性组合(比如对于未定元 a 和 b , 它们是 $a * b$), 假定 q 次乘法运算就足以计算 P 中的所有表达式。证明存在一个包含 q 次

乘法的过程, 其对 k 个 P 中的可以通过单乘法计算的表达式进行计值。

- * 12.10 证明: 至少需要 3 次乘法操作来完成表达式 ac 和 $bc + ad$ 的计算 [提示: 采用 12.9 的结论]。
- * 12.11 证明: 至少需要 $3k$ 次乘法运算以完成有 $2k$ 个表达式的集合的计算, 其中表达式形如: $a_i c_i$ 和 $b_i c_i + a_i d_i$, 其中 $1 \leq i \leq k$ 。

习题 12.12 ~ 12.14 涉及两个 2×2 矩阵 A 和 B 相乘的问题。希望采用非可交换乘法计算 AB 乘积中的四个表达式:

$$\begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \begin{bmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{bmatrix} = \begin{bmatrix} a_{11}b_{11} + a_{12}b_{21} & a_{11}b_{12} + a_{12}b_{22} \\ a_{21}b_{11} + a_{22}b_{21} & a_{21}b_{12} + a_{22}b_{22} \end{bmatrix}$$

- ** 12.12 证明: 如果存在适用于任意环的、需要 q 次乘法的操作, 则存在另一个计算, 采用 q 次乘法操作, 且适用于整数模 2 的环, 这里所有的乘法采用形式:

$$(a_{i,j_1} + a_{i,j_2} + \cdots + a_{i,j_k}) * (b_{m,n_1} + b_{m,n_2} + \cdots + b_{m,n_l})$$

- * 12.13 证明: 如果乘法的左边是如上道习题中的单个 a_{ij} (例如 $k=1$), 则计算至少需要 7 次乘法操作。[提示: 令 $a_{ij}=0$, 采用习题 12.11 取 $k=2$ 的结果。]
- ** 12.14 证明: 如果左边不存在单个的 a_{ij} , 则还需要 7 次乘法操作来完成。[提示: 证明对每种情况, 都有一个条件的集合, 例如 $a_{ij}=0$ 或者 $a_{ii}=a_{nn}$, 使得其中一个乘法等于 0, 剩下的计算就是解决习题 12.11 针对 $k=2$ 的情况。]
- ** 12.15 证明: 2×2 矩阵乘以 $2 \times n$ 矩阵需要 $\lceil 7n/2 \rceil$ 次乘法操作。
- ** 12.16 给出一个采用 $mn + m + n$ 次乘法操作的 $n \times 2$ 矩阵乘以 $2 \times m$ 矩阵的算法, 假定标量乘法形成交换环。把这个结论与习题 12.15 的技术相结合, 证明矩阵乘法在不可交换情况下需要更多的乘法操作。

定义 令 R 为环, 且 a_1, \dots, a_n 和 b_1, \dots, b_n 为两个不相交的未定元集合。双线性型是如下形式的表达式

$$\sum_{i,j} r_{ij} a_i b_j$$

其中 r_{ij} 为 R 中的元素。

- ** 12.17 a) 证明对于任意双线性型的集合都存在一个计算过程, 其需要最少的乘法次数, 且所有的乘法介于 a 的线性函数及 b 的线性函数之间;
- b) 证明存在最小的关于乘法运算的计算过程, 每个乘法计算都介于 a 的线性函数及 b 的线性函数之间 (注意: 如果把这部分限定在交换环时, 这个结论不成立);
- c) 给定普通的允许除法操作的计算过程, 证明存在一个计算过程其需要最少的乘法操作, 而且不需要使用除法操作。
- ** 12.18 令 R 为不可交换环, 且令 \mathbf{a}, \mathbf{b} 和 \mathbf{x} 为未定元的列向量 $[a_1, \dots, a_m]^T$, $[b_1, \dots, b_n]^T$ 和 $[x_1, \dots, x_p]^T$ 。令 X 为 x_i 的线性和的矩阵。习题 12.17 表明计算双线性型的集合 $(\mathbf{a}^T X)^T$ 可以表示为

$$M(P\mathbf{a} \cdot Q\mathbf{x})$$

其中 M, P 和 Q 为元素取自 F 的矩阵。定义表达式集合 $(\mathbf{a}^T X)^T$ 的左对偶为表达式集合 $(\mathbf{b}^T X^T)^T$, 定义 $M(P\mathbf{a} \cdot Q\mathbf{x})$ 的 P 对偶操作为集合 $P^T(M^T \mathbf{b} \cdot Q\mathbf{x})$ 。证明 $(\mathbf{a}^T X)^T$ 的任意计算的 P 对偶计算 $(\mathbf{a}^T X)^T$ 的左对偶。

- * 12.19 证明非交换环上的 $m \times n$ 矩阵乘以 $n \times p$ 矩阵所需的最少乘法操作次数相同于 $n \times m$ 矩阵乘以 $m \times p$ 矩阵。[提示: 采用习题 12.18 的结论。]

目前为止的习题包括了代数操作数目下界的材料, 后面的练习包括更为通用的性质, 问

题 12.20 ~ 12.23 与矩阵转置相关。对 $1 \leq i, j \leq n$, 令 a_{ij} 为未定元。考虑一个计算模型, 每个变量是未定元的 n 元组, 计算步骤 $a \leftarrow b\theta c$ 赋值给 n 元组的 a 个未定元, 由出现在 b 或 c 的未定元中选择。 n 元组 b 或 c 的值不变。

* 12.20 证明任意 n 元组集合的计算

$$\{(a_{i1}, \dots, a_{in}) \mid 1 \leq i \leq n\}$$

以 $\{(a_{i1}, \dots, a_{in}) \mid 1 \leq i \leq n\}$ 作为输入集合, 至少需要 $n \log n$ 步。[提示: 对任意的 i, j , 令 S_{ij} 为出现在一个包含 a_{ij} 的 n 元组中下标为 j 的未定元的最大数, 令 $f = \sum_{i=1}^n \sum_{j=1}^n \log s_{ij}$, 考虑当计算步骤变化时 f 的变化。]

** 12.21 把计算过程的步骤限制为形式 $a \leftarrow b\theta c$, 其中 θ 由 b 中循环移位选择的 $n/2$ 个未定元及从 c 的余码位置中获得 $n/2$ 个未定元。寻找计算过程在 $O(n \log n)$ 时间内计算 12.20 的练习。

** 12.22 设 G 为一个有向无环图, 其有 n 个已知源顶点(没有输入边的顶点)和 n 个输出顶点(没有输出边的顶点)。令 X 和 Y 分别为源顶点与输出顶点的子集, 令 $G(X, Y)$ 为包含从 X 的顶点到 Y 的顶点的路径上所有有向边形成的子图。 $G(X, Y)$ 的容量为分裂 X, Y 所需移除的最少顶点数(与入边和出边数)。假定任意的 X, Y 和图 $G(X, Y)$ 具有容量 $\min(\|X\|, \|Y\|)$, 证明对于所有的 n , 存在有 $cn \log n$ 条边的图 G , 其中 c 是一固定的常数。

** 12.23 移相网络(shifting network)是有 n 个输出顶点, 编号为 $0 \sim n-1$, n 个输入顶点, 编号为 $0 \sim n-1$ 的有向图, 对任意的 $0 \leq s \leq n-1$, 存在一系列不相交的由 i 源顶点入, 从输出 $i+s$ 顶点模 n 出的路径集合。

a) 证明对任意 n , 存在一个移相网络, 其具有 $2n \log n$ 条边;

b) 证明对移相网络, 渐近地需要 $n \log n$ 条边;

c) 证明可以用移相网络计算矩阵的转置。

定义 令 F 为一个域, x_1, x_2, \dots, x_n 为未定元, 线性计算是如下形式的步骤序列: $a \leftarrow \lambda_1 b + \lambda_2 c$, 这里 a 是变量, b 和 c 是变量或者未定元, λ_1 和 λ_2 为 F 的元素, 称为常量。

** 12.24 令 F 为复数域, 令 A 为元素取自 F 的矩阵且 $\mathbf{x} = [x_1, \dots, x_n]^T$ 。证明 $A\mathbf{x}$ 的任意线性计算需要 $\log[\det(A)]/\log(2c)$ 步, 其中 c 为出现在计算过程中的任意常数的最大模[⊖]。

* 12.25 证明对 $[x_1, \dots, x_n]$ 进行傅里叶变换的线性计算过程, 其常数的模小于等于 1, 需要 $\frac{1}{2} n \log n$ 步。

以下 6 个问题主要考虑实现布尔函数所需的两输入的门的最小数量。

* 12.26 证明任意 n 个变量的布尔函数可以至多使用 $n2^n$ 个两输入的门电路实现。

* 12.27 证明对于任意 n , 存在具有 n 个变量的布尔函数, 使得它的实现至少大约需要 $2^n/n$ 个两输入的门电路。

** 12.28 目前对特定布尔函数的实现复杂度的了解甚少, 但是, 若假设将实现限制在树这样的网络中, 可以证明某些函数需要 $n^2/\log n$ 个门电路。证明以下条件足以说明具有 n 个变量的布尔函数需要 $n^2/\log n$ 个门电路。

条件: 令 $f(x_1, \dots, x_n)$ 为具有 n 个变量的布尔函数。把 x_i 划分成为 $b = n/\log n$ 块, 每块具有 $\log n$ 个变量。赋值(0 或 1)给 $b-1$ 块中的各个变量, 总共有 2^{b-1} 种赋值方式。结果得到具有 $\log n$ 个变量的布尔函数, 其是 2^b 个具有 $\log n$ 个变量的布尔函数之一。所得的实际函数取决于给 $b-1$ 块变量所赋的值。令 B_i 为设置为非零的块。如果 B_i 中变量的 $2^{b-1}/n$

⊖ $a+bi$ 的模是 $\sqrt{a^2+b^2}$ 。

种不同函数可以通过对其他变量赋予适当的0和1而得。还要假设对于任意的 i 上述陈述为真, $1 \leq i \leq \log n$, 则任意针对 f 的树形网络需要 $n^2/\log n$ 个门电路。

- ** 12.29 令 $m = 2 + \log n$ 。令 $\{t_{ij} \mid 1 \leq i \leq n/m, 1 \leq j \leq m\}$ 为0-1值的 m 维向量的集合, 其中每个 t_{ij} 至少有两个成分的值是1。令

$$f(x_{11}, x_{12}, \dots, x_{n/m, m}) = \bigoplus_{\substack{1 \leq i \leq n/m \\ 1 \leq j \leq m}} x_{ij} \cdot \left[\bigoplus_{\substack{1 \leq k \leq n/2 \\ k \neq i}} \prod_{\substack{1 \leq l \leq m \\ \text{such that } t_{il} = 1}} x_{kl} \right]$$

其中 t_{il} 为 t_i 的第 l 个成分, 证明采用树形网络实现 f 至少需要 $n^2/\log n$ 个门电路。

- * 12.30 构造其他的布尔函数, 若以树形网络实现它们分别需要(a) $n^{3/2}$ 个门电路和(b) $n^2/\log n$ 个门电路。

- 12.31 令 $S_i(x_1, \dots, x_n)$ 为对称布尔函数, 当且仅当 i 个 x_j 的值为1, 它的值才是1。

a) 给出一种有尽可能少的两输入门的 $S_3(x_1, \dots, x_n)$ 的实现方法;

b) 寻找一种 $S_3(x_1, \dots, x_n)$ 的树形实现, 其采用尽可能少的两输入门。

- ** 12.32 例 12.9 已经证明了计算任意的 n 次多项式需要 n 次乘法操作。试给出一个只需要 n 次乘法的系数为实数的 n 个变量的特别多项式。

- ** 12.33 证明对于由正整数构成的含有 MIN 操作和 + 操作的半环上的矩阵乘法需要 cn^3 次操作, 其中常数 c 满足 $0 < c < 1$ 。

- ** 12.34 证明基于 AND 和 OR 运算的布尔矩阵的乘法需要 n^3 次操作。

习题 12.35 和 12.36 主要考虑多项式的表示, 使多项式的计值大约能在 $n/2$ 次乘法操作内完成。

- ** 12.35 令 $p(x)$ 为 n 次多项式, 其中 n 是奇数且大于 1, 表示为以下形式

$$p(x) = (x^2 - \alpha)q(x) + bx + c$$

a) 证明存在合适的 α , 使得 $b = 0$;

b) 证明存在合适的 α_i 和 β_i , 使得 $p(x)$ 可以表示为

$$p(x) = (x^2 - \alpha_1)[(x^2 - \alpha_2)[\dots] + \beta_2] + \beta_1$$

因此, 用 α_i 和 β_i 表示的 $p(x)$ 的计算可以在 $\lfloor n/2 \rfloor + 2$ 次乘法操作内完成;

c) 在(b)中 α_i 可能是复数, 这种困难应该如何处理? [提示: 选取某个合适的 c , 用 $y + c$ 替换 $p(x)$ 中的 x , 不在 $x = x_0$ 处计算 $p(x)$, 而是在 $y = x_0 - c$ 处计算多项式 $p'(y)$ 。]

习题 12.35 并没有提供方法计算 α_i , 另外, α_i 有可能不是有理数。下面的问题主要研究如何克服这些困难。

- ** 12.36 令 m 为整数, 对 $1 \leq l \leq m$, 令

$$r_l = \frac{(m-l)^2 + (l-1)^2 + 3m + 1}{2}$$

同时令 $r_j = m - l + j + 1, 1 \leq j \leq l$ 。以 α 和 β 为参数定义一个计算链, 形式为

$$c_{11}(x) \leftarrow (x'^n + \alpha_{11})(x'^n + \beta_{11})$$

$$c_{12}(x) \leftarrow (c_{11} + \alpha_{12})(x'^n + \beta_{12})$$

⋮

$$c_{ll}(x) \leftarrow (c_{l,l-1} + \alpha_{ll})(x'^n + \beta_{ll})$$

令 $p(x) = \sum_{i=1}^l c_{ii}(x)$ 。

a) 证明 x 最高次数的系数取决于 $\sum_{i=j}^l r_{ii}$ 给出的 α_{ij} , 同时 x 最高次数的系数也取决于 $\sum_{i=0}^l r_{ii} - r_{ij}$ 给出的 β_{ij} ;

b) 证明所有的多项式次数小于等于 $m(m+1) + 1$, 也就是领先的系数可以被 $m(m+1) + 1$ 个参数表示, 其满足(i)多项式可以在 $m(m+1)/2 + O(m)$ 次乘法内从参数计算出来;
(ii)参数是系数的有理函数。

研究性问题

对于线性状结构的程序中,可使用有向无环图与其关联。图的顶点表示输入和变量。假设 $f \leftarrow g * h$ 为一个步骤,则存在一条从 g 到 f 及 h 到 f 的有向边。观察程序执行的步骤数恰好等于边数的两倍。因此,从程序得到的有向无环图的边数的下界也给出了问题线性状程序执行步骤的下界。

- 12.37 对于大的 n ,证明每个具有 n 个源顶点和 n 个输出顶点、满足习题 12.22 容量条件的有向无回路图至少有 $n \log n$ 条边,或者给出反例。
- 12.38 由两个 n 位整数(假定采用按位计算)相乘的线性状程序所产生的每个图满足习题 12.22 中的容量条件吗?

文献和注释

定理 12.2 及本节所描述问题的通用公式方法取自 Winograd[1970a]。定理 12.1 和 12.3 取自 Fiduccia[1971]。如果没有除法, n 次多项式的计算需要 n 次乘法操作,该结论来自(Knuth)[1969]A. Garsia。Pan[1966]把该结论扩展到乘法运算。Motzkin[1955]证明了经预处理后多项式的计算需要 $\lceil n/2 \rceil + 1$ 次乘法操作。Ostrowski[1954]在这个方向做了一些初步的研究。Winograd[1970b]证明了对于复数乘法至少需要三次乘法操作。

习题 12.7 由 Winograd 提出[1973]。关于矩阵乘法问题下界限的材料(习题 12.9 ~ 12.16)取自 Hopcroft 和 Kerr[1971]。习题 12.17 的结果来自多个人的独立研究。其中 c 部分归功于 S. Winograd 和 P. Ungar 的私人信件。习题 12.18 和 12.19 取自 Hopcroft 和 Muszinski[1973]。习题 12.20 ~ 12.22 及习题 12.37 和 12.38 则基于 R. Floyd 的讨论。习题 12.24 及 12.25 取自 Morgenstern[1973],习题 12.28 和 12.29 取自 Nečiporuk[1966],习题 12.30(a)取自 Harper 和 Savage[1972]。习题 12.32 取自 Strassen[1974]。习题 12.33 取自 Kerr[1970]。习题 12.34 取自 Pratt[1974]。习题 12.35 取自 Eve[1964]。习题 12.36 取自 Rabin 和 Winograd[1971]。

关于求解算术运算下界的其他文献可参阅 Borodin 和 Cook[1974]和 Kedem[1974]。

附录 A 算法的 C/C++ 代码

本书使用一种称为 Pidgin ALGOL 的语言来描述算法。Pidgin 意思是“由两种或多种语言混合而成的一种简化的说话方式，语法和词汇较初等，用于不同语言者进行交流”。算法的开发可分为设计、分析和实现三个环节，该过程可能不断重复，直到需求满足。一般说来，算法的描述只要清晰明了就可以了。书中所用的 Pidgin ALGOL 语言对于算法描述很充分，但 ALGOL 是一个比较古老的语言，如果读者想测试书中给出的例子，要找到一个 ALGOL 语言的编译器是一件不容易的事（译者推荐 PLT Scheme 带的 ALGOL 语言解释器，读者可参阅本书译者的另一本译作《程序设计方法》，人民邮电出版社 2003）。因此，在翻译本书的时候，经和机械工业出版社协商，决定本书的代码保持原貌，另将书中的算法以 C/C++ 实现作为译本的附录，希望对读者有所帮助。

本附录包括了本书近 40 个算法的 C(C++) 代码。

读者可依据 GNU 通用公共授权条款规定，发布与/或修改本代码；代码系基于教学目的而发布，不负任何担保责任；亦无对适售性或特定目的适用性所为的默示性担保。详情请参照 GNU 通用公共授权。

代码编译运行环境：Windows 98/XP，DevC++ 或其他支持 ISO C/C++ 标准的编译器；Linux/Unix，GCC 或其他支持 ISO C/C++ 标准的编译器。

说明：本附录所有代码都在 Dev-C++ Version 4.9.9.2 下测试通过。Dev-C++ 是一个免费的，几乎完全支持 ISO C/C++ 标准的 C/C++ 程序集成开发环境。下载地址为 <http://www.bloodshed.net/devcpp.html>。

使用方法和测试用例见程序注释。

A.1 计算 n^n 的程序

```
#include <iostream>
#include <stdlib.h>

using namespace std;
// 功能：计算  $n$  的  $n$  次方
// 输入：  $n$ 
// 输出：  $n$  的  $n$  次方
// 测试数据：
// 输入 5
// 测试结果：
// 输出  $5^5 = 3125$ 
int nPowerN()
{
    int r1;
    cin >> r1;
    if(r1 <= 0)
    {
        cout << 0;
    }
    else
    {
        int r2 = r1;
        int r3 = r1 - 1;
        while(r3 > 0)
```

```
        {
            r2 = r2 * r1;
            r3 = r3 - 1;
        }
        cout << r2;
    }
    return 0;
}
int main()
{
    nPowerN();
    system("Pause");
    return 0;
}
```

A.2 识别具有相同数目的1和2的字符串的程序

```
#include <iostream>
#include <stdlib.h>

using namespace std;
// 检查一个输入数字流中的1和2的数目是不是相等
// 数字流以0为结束
// 输入:
// Console 输入一个只包含1和2的数字流, 数字流的输入以0结束
// 输出:
// return bool 值, 如果数字流中1的数目和2的数目相等, 返回 true
// 否则, 返回 false
// 测试数据:
// 112221222211110
// 测试结果:
// true;
bool checkEqual()
{
    int x=0;
    int d=0;
    cin >> x;
    while(x!=0)
    {
        if(x!=1)
        {
            d--;
        }
        else
        {
            d++;
        }
        cin >> x;
    }
    if(d==0)
    {
        cout << 1;
        return true;
    }
    return false;
}
int main()
{
    checkEqual();
    system("pause");
    return 0;
}
```

A.3 中序遍历的递归和非递归过程

```

#include <iostream>
#include <stack>

using namespace std;
const int EmptyPosition = -1;
//原书图 2-9、图 2-10,
// 树的中序遍历(运用递归)
// 输入:
//     vertex: 开始中序遍历的顶点
//     leftson: 数组, 用来记录每个顶点的左儿子顶点
//     rightson: 数组, 用来记录每个顶点的右儿子顶点
//     number: 数组, 用来记录各个顶点中序遍历到的序号
// 测试函数:
//     void testInorder()
// 测试数据:
//     输入以下树:
//
//           1
//          / \
//         2   6
//        / \ / \
//       3  4 7  8
//        \      \
//         5      9
// 测试结果:
//     中序遍历结果:
//     3 2 4 5 1 7 6 8 9
void inorder(int vertex, int* leftson, int* rightson, int* number, int countNode)
{
    static int count = 1;
    if (vertex >= countNode || vertex < 0)
    {
        return;
    }
    if (leftson[vertex] != EmptyPosition)
    {
        inorder(leftson[vertex], leftson, rightson, number, countNode);
    }
    number[vertex] = count;
    cout << vertex << " \t";
    count++;
    if (rightson[vertex] != EmptyPosition)
    {
        inorder(rightson[vertex], leftson, rightson, number, countNode);
    }
}
// 不使用递归的树的中序遍历算法实现。
// 输入:
//     root: 开始中序遍历的树顶点。
//     leftson: 数组, 记录各个顶点的左儿子顶点。
//     rightson: 数组, 记录各个顶点的右儿子顶点。
//     number: 记录各个顶点在被中序遍历时的次序序号。
// 测试函数:
//     void testInorder()
// 测试数据:
//     输入以下树:
//
//           1
//          / \

```

```

//          2      6
//        /  \    / \
//       3    4    7   8
//          \    \
//         5      9
// 测试结果:
//    中序遍历结果:
//    3 2 4 5 1 7 6 8 9
void inorderNonRecursive(int* leftson, int* rightson, int* number, int root)
{
    static int count = 1;
    int vertex = root;
    stack<int> s;

    while(1)
    {
        while(vertex != EmptyPosition)
        {
            s.push(vertex);
            vertex = leftson[vertex];
        }

        if(! s.empty())
        {
            vertex = s.top();
            s.pop();
            number[vertex] = count;
            cout << vertex << " \t";
            count ++;
            vertex = rightson[vertex];
        }
        else
        {
            break;
        }
    }
}

void testInorder()
{
    int number[10];
    memset(number, 0, sizeof(int) * 10);
    int leftson[10];
    int rightson[10];

    //build the tree;
    leftson[0] = rightson[0] = EmptyPosition;
    leftson[1] = 2;
    rightson[1] = 6;
    leftson[2] = 3;
    rightson[2] = 4;
    leftson[3] = rightson[3] = EmptyPosition;
    leftson[4] = EmptyPosition;
    rightson[4] = 5;
    leftson[5] = rightson[5] = EmptyPosition;
    leftson[6] = 7;
    rightson[6] = 8;
    leftson[7] = rightson[7] = EmptyPosition;
    leftson[8] = EmptyPosition;
    rightson[8] = 9;
    leftson[9] = rightson[9] = EmptyPosition;

    //inorder with recursive.

```

```

    inorder(1, leftson, rightson, number, 10);
    cout << endl;
    memset(number, 0, sizeof(int)* 10);
    // inorder without recursive
    inorderNonRecursive(leftson, rightson, number, 1);
    cout << endl;
}
int main()
{
    testInorder();
    system("pause");
    return 0;
}

```

A.4 查找 MAX 和 MIN 元素的算法

```

#include <iostream>
using namespace std;
// 参考 2.6 节, 查找一个数组中的最大元素
// 输入:
//   intarray: 输入数组
//   sizeofArray: 输入数组的元素个数
// 输出:
//   返回找到的最大元素
// 测试数据:
//   {1, 3, 5, 4, 2, 8, 9, 1};
// 测试结果:
//   找到最大元素为 9
int findMax(int* intarray, int sizeofArray)
{
    if(intarray == NULL || sizeofArray == 0)
    {
        return 0;
    }
    int max = intarray[0];
    for(int i = 0; i < sizeofArray; i++)
    {
        if(intarray[i] > max)
        {
            max = intarray[i];
        }
    }
    return max;
}
struct MaxMin
{
    int max;
    int min;
    MaxMin()
    {
        max = 0;
        min = 0;
    }
    MaxMin(int _max, int _min)
    {
        max = _max;
        min = _min;
    }
};
// 找到 MAX 和 MIN 过程

```

```

// 找出一个集合中的最大元素和最小元素
// 输入:
//   set: 用数组表示的集合
//   size: 集合中的元素个数
// 输出:
//   返回一个结构体的指针, 结构体中包括了集合的最小元素和最大元素
// 测试函数:
//   void testFindMaxMin()
// 测试数据:
//   {1, 3, 5, 4, 2, 8, 9, 1};
// 测试结果:
//   (max, min) = (9, 1);
MaxMin* findMaxMin(int* set, int size)
{
    MaxMin * result = new MaxMin();
    if(size == 2)
    {
        result->max = max(set[0], set[1]);
        result->min = min(set[0], set[1]);
    }
    else
    {
        int* subset1 = set;
        int* subset2 = set + size/2;
        MaxMin* subresult1 = findMaxMin(subset1, size/2);
        MaxMin* subresult2 = findMaxMin(subset2, size - size/2);
        result->max = max(subresult1->max, subresult2->max);
        result->min = min(subresult1->min, subresult2->min);
    }
    return result;
}

void testFindMaxMin()
{
    int a[8] = {1, 3, 5, 4, 2, 8, 9, 1};
    cout << "max = " << findMax(a, 8) << endl;
    MaxMin * temp = findMaxMin(a, 8);
    cout << "(max, min) = (" << temp->max << ", " << temp->min << ")" << endl;
}

int main()
{
    testFindMaxMin();
    system("pause");
    return 0;
}

```

A.5 归并排序算法

```

#include <iostream>
using namespace std;
//归并两个排序的数组
//输入:
//   a: 数组, 已升序排序
//   sizeofa: 数组 a 的元素个数
//   b: 数组, 已升序排序
//   sizeofb: 数组 b 的元素个数
//输出:
//   返回归并后的数组, 归并后的数组仍然保持升序排序
int* merge(int* a, int sizeofa, int* b, int sizeofb)
{
    int* mergeresult = new int[sizeofa + sizeofb];

```

```

int k=0;
int i=0;
int j=0;
for(; i<sizeofa && j<sizeofb; k++)
{
    if(a[i] <= b[j])
    {
        mergeresult[k] = a[i];
        i++;
    }
    else
    {
        mergeresult[k] = b[j];
        j++;
    }
}
if(i<sizeofa)
{
    for(; i<sizeofa; i++)
    {
        mergeresult[k] = a[i];
        k++;
    }
}
if(j<sizeofb)
{
    for(; j<sizeofb; j++)
    {
        mergeresult[k] = b[j];
        k++;
    }
}
if(a! = NULL)
{
    delete[]a;
}
if(b! = NULL)
{
    delete[]b;
}
return mergeresult;
}
//归并排序的算法实现:
//输入:
//    x: 待排序的数组
//    i: 排序范围的开始位置
//    j: 排序范围的结束位置
//输出:
//    返回排序后的数组 (升序排序)
//测试函数:
//    void testMergeSort()
//测试数据:
//    {2, 34, 5, 3, 65, 7, 88, 2, 8, 36, 92};
//测试结果:
//    {2, 2, 3, 5, 7, 8, 34, 36, 65, 88, 92};
int* mergesort(int* x, int i, int j)
{
    if(i == j)
    {
        int* result = new int[1];

```



```

        result[0] = x[i];
        return result;
    }
    else
    {
        int m = i + j - 1;
        m /= 2;
        return merge(mergesort(x, i, m), m - i + 1, mergesort(x, m + 1, j), j - m);
    }
}

void testMergeSort()
{
    int a[] = {2, 34, 5, 3, 65, 7, 88, 2, 8, 36, 92};
    int* result = mergesort(a, 0, 10);
    for(int i = 0; i < 11; i++)
    {
        cout << result[i] << " \t";
    }
    cout << endl;
    if(result != NULL)
    {
        delete[] result;
    }
}

int main()
{
    testMergeSort();
    system("pause");
    return 0;
}

```

A.6 对矩阵乘法排序的动态规划算法

```

#include <iostream>
using namespace std;
// 计算矩阵相乘的最小计算开销
// 输入:
//     m: 二维数组  $(n+1) \times (n+1)$ , 初始化为全 0;
//     r: 数组, 空间大小为  $n+1$ , 矩阵  $M_i$  的维数维  $r[i-1] * r[i]$ ;
//     n: 共有  $n$  个矩阵相乘
// 输出:
//     m:  $m[0][*]$  与  $m[*][0]$  空出不用,
//     计算后  $m[i][j]$  代表矩阵  $M_i \times \dots \times M_j$  的最小运算开销
// 测试函数:
//     void testCalMinCost()
// 测试数据:
//     r = {10, 20, 50, 1, 100};
//     矩阵 1:  $10 * 20$ 
//     矩阵 2:  $20 * 50$ 
//     矩阵 3:  $50 * 1$ 
//     矩阵 4:  $1 * 100$ 
// 测试结果:
//     矩阵 m:
//         0      10000   1200    2200
//         0       0     1000    3000
//         0       0       0     5000
//         0       0       0       0
void calMinCost(int* * m, int* r, int n)
{
    for(int i = 1; i <= n; i++)

```

```

{
    m[i][i] = 0;
}
for(int l = 1; l <= n-1; l++)
{
    for(int i = 1; i <= n-l; i++)
    {
        int j = i+1;
        m[i][j] = m[i][i] + m[i+1][j] + r[i-1] * r[i] * r[j];
        for(int k = i; k < j; k++)
        {
            if(m[i][j] > m[i][k] + m[k+1][j] + r[i-1] * r[k] * r[j])
            {
                m[i][j] = m[i][k] + m[k+1][j] + r[i-1] * r[k] * r[j];
            }
        }
    }
}
}

void testCalMinCost()
{
    int r[] = {10, 20, 50, 1, 100};
    int ** m = new int * [5];
    for(int i = 0; i < 5; i++)
    {
        m[i] = new int[5];
        memset(m[i], 0, sizeof(int) * 5);
    }
    calMinCost(m, r, 4);
    for(int i = 1; i < 5; i++)
    {
        for(int j = 1; j < 5; j++)
        {
            cout << m[i][j] << '\t';
        }
        cout << endl;
    }
}

int main()
{
    testCalMinCost();
    system("pause");
    return 0;
}

```

A.7 定长串的字典排序算法

```

#include <iostream>
#include <queue>
using namespace std;
// Fig 3.1 Lexicographic sort algorithm
typedef int * pInt;
// k is the number of components in one tuple, n is the number of the tuples,
// 0 to m-1 is the range of the components.
typedef queue<pInt> Bucket;
// 定长的字典序排序
// 输入:
// A: 待排序的整数序列的数组。每个元素是一个整数序列, 按照这些整数的字典序排序
// n: 整数序列的个数
// k: 整数序列的长度

```

```

//      m: 整数序列中的每个整数范围从 0 到 m-1。
// 输出:
//      计算后的 A 是已经排序好的整数序列。
// 测试函数:
//      void testLexicographicSort();
// 测试数据:
//      {1, 2, 3, 5, 6};
//      {2, 3, 6, 9, 2};
//      {1, 3, 5, 4, 3};
//      {2, 3, 4, 5, 6};
//      {1, 3, 4, 8, 4};
// 测试结果:
//      排序后的结果:
//      {1      2      3      5      6}
//      {1      3      4      8      4}
//      {1      3      5      4      3}
//      {2      3      4      5      6}
//      {2      3      6      9      2}
void LexicographicSort(pInt * A, int n, int k, int m)
{
    Bucket * buckets = new Bucket[m];
    queue<pInt> q;
    for(int i=0; i<n; i++)
    {
        q.push(A[i]);
    }
    for(int j=k-1; j>=0; j--)
    {
        while(! q.empty())
        {
            pInt a=q.front();
            q.pop();
            buckets[a[j]].push(a);
        }
        for(int l=0; l<m; l++)
        {
            while(! buckets[l].empty())
            {
                q.push(buckets[l].front());
                buckets[l].pop();
            }
        }
    }
    if(buckets != NULL)
    {
        delete []buckets;
    }
    // showing the result:
    while(! q.empty())
    {
        for(int i=0; i<k; i++)
        {
            cout<<q.front()[i]<<'\\t';
        }
        cout<<endl;
        q.pop();
    }
}
void testLexicographicSort()
{

```

```

pInt * a=new pInt[5];
int a1[]={1, 2, 3, 5, 6};
a[0]=a1;
int a2[]={2, 3, 6, 9, 2};
a[1]=a2;
int a3[]={1, 3, 5, 4, 3};
a[2]=a3;
int a4[]={2, 3, 4, 5, 6};
a[3]=a4;
int a5[]={1, 3, 4, 8, 4};
a[4]=a5;
LexicographicSort(a, 5, 5, 10);
}
int main()
{
    testLexicographicSort();
    system("pause");
    return 0;
}

```

A.8 不同长度的串的字典排序算法

```

#include <iostream>
#include <queue>
#include <vector>
using namespace std;
// Fig 3.2 Lexicographic sort of strings of varing length
typedef string * pString;
typedef queue <pString> StrBucket;
// 变长的字符串的字典序排序
// 输入:
//   strings: 由字符串组成的数组, 待排序
//   strCnt: strings 中的字符串的个数
//   m: 每个字符的范围在 0~m 之间
// 输出:
//   计算之后得到的 strings 就是排序后的字符串数组
// 测试函数:
//   void testVaryLenLexicographicSort();
// 测试数据:
//   "abc", "kdjkelg", "abdkdj", "dkehgke", "hello", "heloo", "dkekk";
// 测试结果:
//   字典序排序后的结果:
//   "abc", "abdkdj", "dkehgke", "dkekk", "hello", "heloo", "kdjkelg";
void VaryLenLexicographicSort(string * strings, int strCnt, int m=256)
{
    unsigned int maxLength=0;
    for(int i=0; i<strCnt; i++)
    {
        if(maxLength<strings[i].length())
        {
            maxLength=strings[i].length();
        }
    }
    vector <pString> * length=new vector <pString> [maxLength+1];
    vector <char> * nonempty=new vector <char> [maxLength];
    for(int i=0; i<strCnt; i++)
    {
        (length[strings[i].length()]).push_back(&(strings[i]));
        for(int j=0; j<strings[i].length(); j++)
        {

```

```

        nonempty[j].push_back(strings[i].at(j));
    }
}
queue<pString> q;
StrBucket * buckets = new StrBucket[m];

for(int l=maxlength-1; l >= 0; l--)
{
    for(int x=0; x<length[l+1].size(); x++)
    {
        q.push(length[l+1].at(x));
    }
    while(! q.empty())
    {
        pString ps = q.front();
        buckets[ps->at(l)].push(ps);
        q.pop();
    }
    sort(nonempty[l].begin(), nonempty[l].end());
    for(int j=0; j<nonempty[l].size(); j++)
    {
        int val = nonempty[l].at(j);
        while(! buckets[val].empty())
        {
            pString ps = buckets[val].front();
            q.push(ps);
            buckets[val].pop();
        }
    }
}
while(! q.empty())
{
    cout << * (q.front()) << endl;
    q.pop();
}
if(length != NULL)
{
    delete [] length;
}
if(nonempty != NULL)
{
    delete [] nonempty;
}
if(buckets != NULL)
{
    delete [] buckets;
}
}

void testVaryLenLexicographicSort()
{
    string * strings = new string[7];
    strings[0] = "abc";
    strings[1] = "kdjkelg";
    strings[2] = "abdkdj";
    strings[3] = "dkehgke";
    strings[4] = "hello";
    strings[5] = "heloo";
    strings[6] = "dkekk";
    VaryLenLexicographicSort(strings, 7);
}

```

```

int main()
{
    testVaryLenLexicographicSort();
    system("pause");
    return 0;
}

```

A.9 快速排序程序

```

#include <iostream>
using namespace std;
// Fig 3.8 Partitioning S into S1 and S2 union S3, in place.
// 划分数组，以 s[0] 为界，小于 s[0] 的置于数组左段，大于等于 s[0] 的元素置于数组右段
// 输入：
//     s: 待划分的数组
//     f: 划分范围的开始下标
//     l: 划分范围的结束下标
// 输出：
//     n1: 小于 s[0] 的元素个数
//     s: 划分好的数组
void partitionArray(int * s, int f, int l, int &n1)
{
    n1 = 0;
    int pivot = f;
    int i = f + 1;
    int j = l;
    while(i <= j)
    {
        while(j >= f && s[j] >= s[pivot])
        {
            j--;
        }
        while(i <= l && s[i] < s[pivot])
        {
            i++;
        }
        if(j >= f && i <= l && i < j)
        {
            // 交换 s[i] 和 s[j]
            int t = s[i];
            s[i] = s[j];
            s[j] = t;
            i++;
            j--;
        }
    }
    if(j >= f)
    {
        // 交换 s[pivot] 和 s[j]
        int t = s[pivot];
        s[pivot] = s[j];
        s[j] = t;
        n1 = j - f;
    }
    else
    {
        n1 = 0;
    }
}
// Fig. 3.7 Quick sort program

```

```

//快速排序算法的递归部分
//输入:
//    s: 待排序的数组
//    f: 递归覆盖范围的开始下标
//    l: 递归覆盖范围的结束下标
//输出:
//    s: 递归快速排序
void quickSortDriver(int * s, int f, int l)
{
    if(f >= l)
    {
        return;
    }
    else
    {
        int n1 = 0;
        partitionArray(s, f, l, n1);
        quickSortDriver(s, f, f + n1 - 1);
        quickSortDriver(s, f + n1 + 1, l);
    }
}
// 快速排序算法实现
// 输入:
//    s: 待排序的数组
//    n: 数组的大小
// 输出:
//    s: 排序的数组
// 测试函数:
//    void testQuickSort();
// 测试数据:
//    {23, 2, 5, 63, 8, 2, 37, 9, 12, 34, 3};
// 测试结果:
//    {2, 2, 3, 5, 8, 9, 12, 23, 34, 37, 63}
void quickSort(int * s, int n)
{
    quickSortDriver(s, 0, n - 1);
}
void testQuickSort()
{
    int a[] = { 23, 2, 5, 63, 8, 2, 37, 9, 12, 34, 3};
    quickSort(a, 11);
    for(int i = 0; i < 11; i++)
    {
        cout << a[i] << '\t';
    }
    cout << endl;
}
int main()
{
    testQuickSort();
    system("pause");
    return 0;
}

```

A. 10 选择第 k 小元素的算法

```

#include <iostream>
#include <math.h>
using namespace std;
// Fig 3.8 Partitioning S into S1 and S2 union S3, in place.

```

```

// 划分数组, 以 s[0] 为界, 小于 s[0] 的置于数组左段, 大于等于 s[0] 的元素置于数组右段
// 输入:
//   s: 待划分的数组
//   f: 划分范围的开始下标
//   l: 划分范围的结束下标
// 输出:
//   n1: 小于 s[0] 的元素个数
//   s: 划分好的数组
void partitionArray(int * s, int f, int l, int & n1)
{
    n1 = 0;
    int pivot = f;
    int i = f + 1;
    int j = l;
    while(i <= j)
    {
        while(j >= f && s[j] >= s[pivot])
        {
            j--;
        }
        while(i <= l && s[i] < s[pivot])
        {
            i++;
        }
        if(j >= f && i <= l && i < j)
        {
            // 交换 s[i] 和 s[j]
            int t = s[i];
            s[i] = s[j];
            s[j] = t;
            i++;
            j--;
        }
    }
    if(j >= f)
    {
        // 交换 s[pivot] 和 s[j]
        int t = s[pivot];
        s[pivot] = s[j];
        s[j] = t;
        n1 = j - f;
    }
    else
    {
        n1 = 0;
    }
}

// Fig. 3.7 Quick sort program
// 快速排序算法的递归部分
// 输入:
//   s: 待排序的数组
//   f: 递归覆盖范围的开始下标
//   l: 递归覆盖范围的结束下标
// 输出:
//   s: 递归快速排序
void quickSortDriver(int * s, int f, int l)
{
    if(f >= l)
    {
        return;
    }
}

```



```

    }
    else
    {
        int n1 = 0;
        partitionArray(s, f, l, n1);
        quickSortDriver(s, f, f + n1 - 1);
        quickSortDriver(s, f + n1 + 1, l);
    }
}
// 快速排序算法实现。
// 输入:
//   s:   待排序的数组
//   n:   数组的大小
// 输出:
//   s:   排序的数组
// 测试函数:
//   void testQuickSort();
// 测试数据:
//   {23, 2, 5, 63, 8, 2, 37, 9, 12, 34, 3};
// 测试结果:
//   {2, 2, 3, 5, 8, 9, 12, 23, 34, 37, 63}
void quickSort(int * s, int n)
{
    quickSortDriver(s, 0, n - 1);
}
// fig 3.10 Algorithm to select kth smallest element.
// 从数组 s 中找出第 k 大的元素
// 输入:
//   k:   指定需从数组中找出第 k 大的元素
//   s:   容纳所有元素的数组
//   n:   数组中元素的个数
// 输出:
//   返回找到的数组中第 k 大个元素
// 测试函数:
//   void testSelect();
// 测试数据:
//   { 23, 2, 5, 63, 8, 2, 37, 9, 12, 34, 3};
// 测试结果:
//   k=10 从小到大排序位置为 10 的数字是: 37
int select(int k, int * s, int n)
{
    if(k > n)
    {
        cout << "k 不在正确选择范围内" << endl;
        return -1;
    }
    if(n < 50)
    {
        quickSort(s, n);
        return s[k - 1];
    }
    else
    {
        int count = n / 5;

        // 分组处理, 每组最多 5 个元素。
        for(int i = 0; i < count; i++)
        {
            quickSort(s + i * 5, 5);
        }
    }
}

```

```

int * M=new int[count];
int j=2;
for(int i=0; i<count; i++)
{
    M[i]=s[j];
    j+=5;
}
int m=select(static_cast<int>(ceil(count/2.0)), M, count);
int x=0;
int y=n-1;
int * ss=new int[n];
for(int i=0; i<n; i++)
{
    if(s[i]<m)
    {
        ss[x]=s[i];
        x++;
    }
    else
    {
        ss[y]=s[i];
        y--;
    }
}
int nl=x;
if(nl>=k)
{
    return select(k, ss, nl);
}
else
{
    return select(k-nl, ss+nl, n-nl);
}
if(ss!=NULL)
{
    delete [] ss;
}
}
}
// Fig 3.12 Selection algorithm
// 从数组 s 中找出第 k 大的元素
// 输入:
//     k: 指定需从数组中找出第 k 大的元素
//     s: 容纳所有元素的数组
//     n: 数组中元素的个数
// 输出:
//     返回找到的数组中第 k 大个元素。
// 测试函数:
//     void testSelect();
// 测试数据:
//     { 23, 2, 5, 63, 8, 2, 37, 9, 12, 34, 3 };
// 测试结果:
//     k=10 从小到大排序位置为 10 的数字是: 37
int selectRandom(int k, int * s, int n)
{
    if(k>n)
    {
        cout << "k 不在正确选择范围内" << endl;
        return -1;
    }
}

```

```

    if(n == 1)
    {
        return s[0];
    }
    else
    {
        int n1 = 0;
        partitionArray(s, 0, n - 1, n1);
        if(n1 >= k)
        {
            return select(k, s, n1);
        }
        else
        {
            return select(k - n1, s + n1, n - n1);
        }
    }
}

void testSelect()
{
    int a[] = { 23, 2, 5, 63, 8, 2, 37, 9, 12, 34, 3 };
    int k = 10;
    cout << select(k, a, 11) << endl;
    cout << selectRandom(k, a, 11) << endl;
}

int main()
{
    testSelect();
    system("pause");
    return 0;
}

```

A. 11 最小代价生成树算法

```

#include <iostream>
#include <list>
#include <vector>
#include <algorithm>
#include <stack>
#include <stdlib.h>
using namespace std;
struct Edge
{
    int start;
    int end;
    float weight;
    Edge(int _start, int _end, float _weight)
    {
        start = _start;
        end = _end;
        weight = _weight;
    }
    bool operator < (const Edge & e) const
    {
        return this->weight < e.weight;
    }
    bool operator > (const Edge & e) const
    {
        return this->weight > e.weight;
    }
}

```

```

bool operator > = (const Edge & e) const
{
    return this->weight > = e.weight;
}
bool operator < = (const Edge & e) const
{
    return this->weight < = e.weight;
}
bool operator == (const Edge & e) const
{
    return (this->weight == e.weight &&
            this->start == e.start &&
            this->end == e.end);
}

};
int find(int v, int * name, int * father);
void unionSet(int * father, int * name, int * count, int * root,
              int i, int j, int k, int setnamerange);

// Fig. 5.2
// 最小生成树的生成
// 输入:
//     V: 原图的点集 (点从 1 到 nV-1)
//     nV: 原图的点的数目
//     E: 原图的边集
// 输出:
//     生成的最小生成树边集为 T, 点集仍然为 V
// 测试函数:
//     testBuildMinimumCostSpanningTree()
// 测试数据:
//     原图由 4 点组成: 0, 1, 2, 3, 4, 5, 6
//
//           20
//         v0 - - - - - v1
//        / \         / \
//       23/   \1   4/   \15
//        /36   \   / 9   \
//       v5 - - - - - v6 - - - - - v2
//        \   /   \   /
//       28 \   /25   \16 /3
//            \ /   17   \ /
//           v4 - - - - - v3
// 测试结果:
//     生成树点集与原图一致: 0, 1, 2, 3, 4, 5, 6
//
//           v0           v1
//          / \         /
//         23/   \1   4/
//          /   \   / 9
//         v5     v6 - - - - - v2
//                  /
//                  /3
//                 17   /
//                v4 - - - - - v3
// void buildMinimumCostSpanningTree(int * V, int nV, vector<Edge> & E, vector<Edge> & T)
// {
//     T.clear();
//     // 初始化一个集合的集合
//     int * father = new int[nV];
//     int * root = new int[nV];
//     int * count = new int[nV];
//     int * name = new int[nV];

```

```

// 每个点都初始化为一个单独的点集
for(int i=0; i<nV; i++)
{
    father[i] = -1;
    count[i] = 1; // 每个集合 i 中都只有 i 这一个元素
    root[i] = i; // 集合 i 的根顶点就是点 i;
    name[i] = i;
}
// 用堆的形式来存放图中所有的边 E;
make_heap(E.begin(), E.end(), greater<Edge>());
vector<Edge>::iterator edgeEndIt = E.end();
while(true)
{
    if(count[root[find(0, name, father)]] >= nV || edgeEndIt == E.begin())
    { // 所有的点都在同一个点集中了, 表示 T 中已经覆盖了所有的点
        break;
    }
    // 从 E 中选出权重最小的一条边 (v, w);
    Edge e = E.at(0);
    // 从 E 中删除 (v, w);
    pop_heap(E.begin(), edgeEndIt, greater<Edge>());
    edgeEndIt--;
    int w1 = find(e.start, name, father);
    int w2 = find(e.end, name, father);
    if(w1 != w2) // v 和 w 在 VS 中的不同点集 w1 和 w2 中
    {
        // 用点集 w1 union w2 代替 VS 中的点集 w1 和 w2;
        unionSet(father, name, count, root, w1, w2, nV);
        // 将边 (v, w) 加入到 T 中;
        T.push_back(e);
    }
}
}

void testBuildMinimumCostSpanningTree()
{
    int V[] = {0, 1, 2, 3, 4, 5, 6};
    vector<Edge> E;
    Edge e01(0, 1, 20);
    Edge e12(1, 2, 15);
    Edge e23(2, 3, 3);
    Edge e34(3, 4, 17);
    Edge e45(4, 5, 28);
    Edge e50(5, 0, 23);
    Edge e06(0, 6, 1);
    Edge e16(1, 6, 4);
    Edge e26(2, 6, 9);
    Edge e36(3, 6, 16);
    Edge e46(4, 6, 25);
    Edge e56(5, 6, 36);
    E.push_back(e01);
    E.push_back(e12);
    E.push_back(e23);
    E.push_back(e34);
    E.push_back(e45);
    E.push_back(e50);
    E.push_back(e06);
    E.push_back(e16);
    E.push_back(e26);
    E.push_back(e36);
    E.push_back(e46);
}

```

```

    E.push_back(e56);
    vector<Edge> tree;
    tree.clear();
    buildMinimumCostSpanningTree(V, 7, E, tree);
    cout << "In the min cost spanning tree built: " << endl;
    for(vector<Edge>:: iterator it = tree.begin(); it != tree.end(); it++)
    {
        cout << "(" << it->start << ", " << it->end << "): " << it->weight << endl;
    }
}

int main()
{
    testBuildMinimumCostSpanningTree();
    system("pause");
    return 0;
}

// Fig 4.18 Executing instruction FIND(i)
// 在集合中查找某个元素所在的集合名(集合用树的结构来表示)
// 输入:
//   v : 查找顶点 v
// name: 数组, 存放每个根顶点对应的集合的名称。name[v]表示根顶点为 v 的集合的名称
// father: 数组, 存放每个顶点的父顶点, 如 father[v]表示顶点 v 的父顶点
// 输出:
//   Vi 所属于的集合名将打印到控制台(console)上, 并返回 Vi 所属于的集合名
// 测试数据: 给定集合:
//           7
//           / \
//          5   2
//           /
//          1
//           /
//          3
//   查找 v = 3;
// 测试结果: 找出根顶点 7, 打印出集合名 3, 并将集合调整为一下结构
//           7
//           / \ \ \
//          5   2 1 3
int find(int v, int * name, int * father)
{
    list<int> l;
    l.clear();
    // 从节点 v 向上追溯, 知道这个集合的根顶点
    while(father[v] != -1)
    {
        l.push_back(v);
        v = father[v];
    }
    // v 此刻已经是根顶点
    // 调整向上追溯的路径上的所有的顶点, 调整它们使得它们全部都直接连在根上, 减少之后搜索的代价
    for(list<int>:: iterator it = l.begin(); it != l.end(); it++)
    {
        father[* it] = v;
    }
    return name[v];
}

// fig 4.19 Executing instruction UNION(i, j, k)
// 集合合并算法实现(集合用树结构来表示)
// 输入:
//   father: 数组, 存放每个节点的父顶点, 如 father[v]表示顶点 v 的父顶点
//   name: 数组, 存放每个根顶点对应的集合的名称。name[v]表示根顶点为 v 的集合的名称

```

```

// count: 数组, 存放每个根顶点代表的集合中的元素的个数
// count[v]表示根顶点是v的集合的元素的个数
// root: 数组, 存放每个集合对应的根顶点。root[i]表示集合i的根顶点v
// i, j: 给外部名称为i和j的集合求并
// k: 求并后的集合外部名称
// setnamerange: 集合名的范围, 只能从0到(setnamerange-1)
// 输出:
// 运行结束后的数组中存放的就是求并之后的结果
// 测试数据:
// 求并之前的集合如下:
// 集合1: 4    集合2: 7    集合3: 10
//      / \      \      /
//      2  5      6      11
//      / \      / \
//      3  9      1  8
// 将集合2 和集合3 求并, 求并后的新集合为 4
// 测试结果:
// 集合1 不变, 集合4    7
//      / \
//      6  10
//      / \ \
//      1  8  11
void unionSet(int * father, int * name, int * count, int * root, int i, int j, int k, int setnamerange)
{
    if(i >= setnamerange || j >= setnamerange || k >= setnamerange ||
        i < 0 || j < 0 || k < 0)
    {
        cout << "the set name is out of the range. " << endl;
        cout << "please use the set name among 0 to " << setnamerange - 1 << endl;
        return;
    }
    if(root[i] == -1 && root[j] == -1)
    {
        cout << "error: the set i and j are both empty! " << endl;
        return;
    }
    int counti = root[i] == -1 ? 0 : count[root[i]];
    int countj = root[j] == -1 ? 0 : count[root[j]];
    if(counti > countj)
    {
        int t = i;
        i = j;
        j = t;
    }
    int large = root[j];
    int small = root[i];
    if(small != -1) // 小的集合不是空集
    {
        father[small] = large;
        name[small] = -1;
        count[small] = 0;
    }
    count[large] = count[large] + count[small];
    name[large] = k;
    root[i] = -1;
    root[j] = -1;
    root[k] = large;
}

```

A. 12 二分搜索算法

```
#include <iostream>
using namespace std;
// Fig 4.3 Binary search algorithm
// 二分搜索算法
// 输入:
//     a: 待搜索的元素
//     A: 供搜索的数组, 数组中的元素是升序排序的
//     f: 数组搜索范围的开始下标
//     l: 数组搜索范围的结束下标
// 输出:
//     返回 bool 变量, 找到返回 true, 未找到返回 false
// 测试函数:
//     void testbinarySearch();
// 测试数据:
//     {1, 3, 6, 7, 43, 45, 98, 345, 455, 599}
// 测试结果:
//     在数据中搜索: 455 返回 true
//     在数据中搜索: 47 返回 false
bool search(int a, int * A, int f, int l)
{
    if(f > l)
    {
        return false;
    }
    else
    {
        if(a == A[(f + l) / 2])
        {
            return true;
        }
        else if(a < A[(f + l) / 2])
        {
            return search(a, A, f, (f + l) / 2 - 1);
        }
        else
        {
            return search(a, A, (f + l) / 2 + 1, l);
        }
    }
}

void testbinarySearch()
{
    int a[] = {1, 3, 6, 7, 43, 45, 98, 345, 455, 599};
    int searchfor = 455;
    if(search(searchfor, a, 0, 9))
    {
        cout << "找到 " << searchfor << endl;
    }
    else
    {
        cout << searchfor << "不在待搜索数组中!" << endl;
    }
    searchfor = 47;
    if(search(searchfor, a, 0, 9))
    {
        cout << "找到 " << searchfor << endl;
    }
    else
```



```

    {
        cout << searchfor << "不在待搜索数组中!" << endl;
    }
}
int main()
{
    testbinarySearch();
    system("pause");
    return 0;
}

```

A. 13 在二叉查找树中搜索特定元素

```

#include <iostream>
using namespace std;
// Fig 4.5 Searching a binary search tree.
// 二叉查找树顶点定义
struct BSTNode
{
    int value; // 节点的值
    BSTNode * left; // 节点的左子顶点
    BSTNode * right; // 节点的右子顶点
    BSTNode(int _value)
    {
        value = _value;
        left = NULL;
        right = NULL;
    }
    BSTNode()
    {
    }
};
// 在二叉查找树中搜索特定元素
// 输入:
//     a: 待搜索的值
//     v: 开始搜索的根顶点
// 输出:
//     返回 bool 值, 如果在树中找到待搜索的值, 返回 true, 否则, 返回 false
// 测试函数:
//     void testBinarySearchTreeSearch();
// 测试数据:
//     在下列树中搜索特定元素
//           5
//          / \
//         2   7
//        / \ / \
//       1  3 6 10
// 测试结果:
//     搜索 8 : 返回 false
//     搜索 7 : 返回 true
bool search(int a, BSTNode * v)
{
    if(a == v->value)
    {
        return true;
    }
    else if(a < v->value)
    {
        if(v->left != NULL)
        {

```

```

        return search(a, v->left);
    }
    else
    {
        return false;
    }
}
else
{
    if(v->right != NULL)
    {
        return search(a, v->right);
    }
    else
    {
        return false;
    }
}
}
void testBinarySearchTreeSearch()
{
    BSTNode * root = new BSTNode(5);
    BSTNode * cur = root;
    cur->left = new BSTNode(2);
    cur->right = new BSTNode(7);
    cur = cur->left;
    cur->left = new BSTNode(1);
    cur->right = new BSTNode(3);
    cur = root->right;
    cur->left = new BSTNode(6);
    cur->left = new BSTNode(10);
    int searchfor = 7;
    if(search(searchfor, root))
    {
        cout << searchfor << "在二叉查找树中" << endl;
    }
    else
    {
        cout << "树中未找到: " << searchfor << endl;
    }

    searchfor = 8;
    if(search(searchfor, root))
    {
        cout << searchfor << "在二叉查找树中" << endl;
    }
    else
    {
        cout << "树中未找到: " << searchfor << endl;
    }
}
int main()
{
    testBinarySearchTreeSearch();
    system("pause");
    return 0;
}

```

A. 14 计算最优子树根的算法

```
#include <iostream>
```

```

using namespace std;
//二叉查找树顶点定义
struct BSTNode
{
    int value; // 顶点的值
    BSTNode * left; // 顶点的左子顶点
    BSTNode * right; // 顶点的右子顶点
    BSTNode(int _ value)
    {
        value = _ value;
        left = NULL;
        right = NULL;
    }
    BSTNode()
    {
    }
};

// Fig 4.9. Algorithm to compute roots of optimal subtrees.
// 计算最优子树的根的算法实现
// p, q 的数组规模是 n+1; root 的矩阵规模是 n × n
// 输入:
//   p: 数组, p[i] 表示查找集合中第 i 个元素的概率
//   q: 数组, q[0] 表示查找 a < A[1] 的概率, q[i] 表示查找 A[i] < a < a[i+1] 的概率
//   q[n] 表示查找 a > A[n] 的概率
//   root: 初始化为全 0
//   n: 集合 A 中元素的个数为 n-1, 表示时占用数组的 1 - (n-1) 下标, 下标 0 保留不用
// 输出:
//   root: 计算后的最优根顶点。root[i][j] 表示子树 a[i+1]...a[j] 的最佳根顶点
// 测试函数:
//   void testOptimalBinaryTree();
void calRootsOfOptimalBinaryTree(float * p, float * q, int ** root, int n)
{
    typedef float * pFloat;

    pFloat * weight = new pFloat[n];
    pFloat * cost = new pFloat[n];
    for(int i=0; i<n; i++)
    {
        weight[i] = new float[n];
        cost[i] = new float[n];
    }
    for(int i=0; i<n; i++)
    {
        for(int j=0; j<n; j++)
        {
            weight[i][j] = 0;
            cost[i][j] = 0;
        }
    }
    for(int i=0; i<n; i++)
    {
        weight[i][i] = q[i];
    }

    for(int l=1; l<n; l++)
    {
        for(int i=0; i<n-l; i++)
        {
            int j = l + i;
            weight[i][j] = weight[i][j-1] + p[j] + q[j];

```

```

        int k=i+1;
        int m=k;
        float minVal = cost[i][k-1] + cost[k][j];
        for(; k <= j; k++)
        {
            if(cost[i][k-1] + cost[k][j] < minVal)
            {
                minVal = cost[i][k-1] + cost[k][j];
                m = k;
            }
        }
        cost[i][j] = weight[i][j] + cost[i][m-1] + cost[m][j];
        root[i][j] = m;
    }
}

// a is the array of all elements in the set.
//根据计算出来的最佳根顶点, 建立最优子树
//输入:
//    i: 表示树中顶点从 a[i+1]开始
//    j: 表示树中顶点范围到 a[j]结束
//    a: 表示集合, 包含所有的集合元素。其中的元素满足 a[i] < a[i+1]
//    n: 表示集合元素下标从 1 到 n-1, 数组的规模为 n
//输出:
//    返回建立子树的根顶点
BSTNode * buildTree(int i, int j, int ** root, int * a, int n)
{
    if(i >= n || j >= n)
    {
        return NULL;
    }
    int m = root[i][j]; // this is the root of the tree.
    BSTNode * rootnode = new BSTNode();
    rootnode->value = a[m];
    rootnode->right = NULL;
    rootnode->left = NULL;
    if(i < m-1)
    {
        rootnode->left = buildTree(i, m-1, root, a, n);
    }
    if(j > m)
    {
        rootnode->right = buildTree(m, j, root, a, n);
    }
    return rootnode;
}

void inorderPrintTree(BSTNode * root)
{
    if(root == NULL)
    {
        return;
    }
    inorderPrintTree(root->left);
    cout << root->value << endl;
    inorderPrintTree(root->right);
}

void testOptimalBinaryTree()
{
    // 初始化 root 矩阵
    int ** root = new int * [5];

```

```

for(int i=0; i<5; i++)
{
    root[i]=new int [5];
    memset(root[i], 0, sizeof(int) * 5);
}
// 初始化 p, q, a;
int a [] = {0, 5, 7, 9, 10};
float p [] = {0, 0.25, 0.125, 0.0625, 0.0625};
float q [] = {0.125, 0.1875, 0.0625, 0.0625, 0.0625};
calRootsOfOptimalBinaryTree(p, q, root, 5);
BSTNode * r=buildTree(0, 4, root, a, 5);
for(int i=0; i<5; i++)
{
    for(int j=0; j<5; j++)
    {
        cout << root[i][j] << '\t';
    }
    cout << endl;
}
inorderPrintTree(r);
}
int main()
{
    testOptimalBinaryTree();
    system("pause");
    return 0;
}

```

A.15 UNION 算法

```

#include <iostream>
using namespace std;
// Fig 4.14 Implementation of UNION instruction
// union the set i and the set j, set k is the destination set.
// 集合求并的算法实现
// 输入:
//   internal_name: internal_name[i] 表示外部名为 i 的集合的内部名
//   external_name: external_name[i] 表示内部名为 i 的集合的外部名
//   size: size[i] 表示内部名称为 i 的集合的元素个数
//   list: list[i] 表示内部名称为 i 的集合的第一个元素在 refer 中的下标
//   refer: refer[i] 中放置的是元素 i 所属于的集合的内部名称
//   next: next[i] 中放置的是元素 i 的下一个元素在 refer 中的下标
//   i, j: 给外部名称为 i 和 j 的集合求并
//   k: 求并后的集合外部名称
// 输出:
//   以上的各个数组中的内容被更新, 表示了最新求并之后的集合
// 测试函数:
//   void testUnionSet();
// 测试数据:
//   集合 1 = {0, 3, 5, 7};
//   集合 2 = {2, 4, 8}
//   集合 3 = {6};
// 测试结果:
//   将集合 1, 2 取并成为集合 4 后
//   集合 3: {6,}
//   集合 4: {2, 4, 8, 0, 3, 5, 7,}
void unionSet(int * internal_name, int * external_name, int * size, int * list, int * refer,
int * next,
                int i, int j, int k)
{

```

```

int a = internal_name[i];
int b = internal_name[j];
int last = -1;
if(size[a] > size[b])
{
    // exchange the role of a and b
    int t = a;
    a = b;
    b = t;
}
int element = list[a];
while(element != -1)
{
    refer[element] = b;
    last = element;
    element = next[element];
}
next[last] = list[b];
list[b] = list[a];
size[b] = size[a] + size[b];
internal_name[k] = b;
external_name[b] = k;
}

void showSets(int * refer, int * next, int * list, int * size, int * external_name, int *
internal_name, int setname)
{
    int inter_name = internal_name[setname];
    cout << "集合" << setname << ": {" ;
    for(int e = list[inter_name]; e != -1; e = next[e])
    {
        cout << e << ", ";
    }
    cout << "}" << endl;
}

void testUnionSet()
{
    int refer [] = {1, -1, 2, 1, 2, 1, 0, 1, 2};
    int next [] = {3, -1, 4, 5, 8, 7, -1, -1, -1};
    int list [] = {6, 0, 2};
    int size [] = {1, 4, 3};
    int external_name [] = {3, 1, 2};
    int internal_name [7] = {-1, 1, 2, 0, -1, -1, -1};
    // 以上数据结构, 表示了集合1 = {0, 3, 5, 7};
    // 集合2 = {2, 4, 8}
    // 集合3 = {6};
    showSets(refer, next, list, size, external_name, internal_name, 1);
    showSets(refer, next, list, size, external_name, internal_name, 2);
    showSets(refer, next, list, size, external_name, internal_name, 3);
    unionSet(internal_name, external_name, size, list, refer, next, 1, 2, 4);
    cout << "将集合1, 2 取并成为集合4后: " << endl;
    showSets(refer, next, list, size, external_name, internal_name, 3);
    showSets(refer, next, list, size, external_name, internal_name, 4);
}

int main()
{
    testUnionSet();
    system("pause");
    return 0;
}

```

A. 16 FIND(i) 的执行

```

#include <iostream>
#include <list>
using namespace std;
// Fig 4.18 Executing instruction FIND(i)
// 在集合中查找某个元素所在的集合名(集合用树的结构来表示)
// 输入:
//     v : 查找顶点 v
// name: 数组, 存放每个根顶点对应的集合的名称。name[v]表示根顶点为 v 的集合的名称
// father: 数组, 存放每个顶点的父顶点, 如 father[v]表示顶点 v 的父顶点
// 输出:
//     Vi 所属于的集合名将被打印到控制台(console)上, 并返回 Vi 所属于的集合名
// 测试函数: testFind()
// 测试数据: 给定集合: (集合名为 3)
//           7
//           / \
//          5   2
//           /
//          1
//           /
//          3
//     查找 v = 3;
// 测试结果: 找出根顶点 7, 打印出集合名 3, 并将集合调整为一下结构
//           7
//           / \ \ \
//          5   2 1 3
int find(int v, int * name, int * father)
{
    list<int> l;
    l.clear();
    // 从节点 v 向上追溯, 知道这个集合的根顶点
    while(father[v] != -1)
    {
        l.push_back(v);
        v = father[v];
    }
    // v 此刻已经是根顶点。
    // 调整向上追溯的路径上的所有的顶点, 调整它们使得它们全部都直接连在根上,
    // 减少之后搜索的代价
    for(list<int>::iterator it = l.begin(); it != l.end(); it++)
    {
        father[* it] = v;
    }
    return name[v];
}
void testFind()
{
    int name[] = {-1, -1, -1, -1, -1, -1, -1, 3};
    int father[] = {-1, 2, 7, 1, -1, 7, -1, -1};
    int vertex = 3;
    int set = find(vertex, name, father);
    cout << "节点" << vertex << " 属于集合: " << set << endl;
}
int main()
{
    testFind();
    system("pause");
    return 0;
}

```

A.17 UNION(i, j, k) 的执行

```

#include <iostream>
#include <list>
using namespace std;
int find(int v, int * name, int * father);
const int EMPTY_ POSITION = -1;
// fig 4.19 Executing instruction UNION(i, j, k)
// 集合求并算法实现(集合用树结构来表示)
// 输入:
//     father: 数组, 存放每个顶点的父顶点, 如 father[v]表示顶点 v 的父顶点
//     name: 数组, 存放每个根顶点对应的集合的名称
//     name[v]表示根顶点为 v 的集合的名称
//     count: 数组, 存放每个根顶点代表的集合中的元素的个数
//     count[v]表示根顶点是 v 的集合的元素的个数
//     root: 数组, 存放每个集合对应的根顶点。root[i]表示集合 i 的根顶点 v
//     i, j: 给外部名称为 i 和 j 的集合求并
//     k: 求并后的集合外部名称
//     setnamerange: 集合名的范围, 只能从 0 到 (setnamerange - 1)
// 输出:
//     运行结束后的数组中存放的就是求并之后的结果
// 测试函数:
//     void testUnionTreeSet()
// 测试数据:
//     求并之前的集合如下:
//     集合 1: 4      集合 2: 7      集合 3: 10
//           / \          \          /
//          2   5          6          11
//           / \          / \
//          3   9          1   8
//     将集合 2 和集合 3 求并, 求并后的新集合为 4
// 测试结果:
//     集合 1 不变, 集合 4      7
//           / \
//          6   10
//           / \ \
//          1   8  11
void unionSet(int * father, int * name, int * count, int * root,
              int i, int j, int k, int setnamerange)
{
    if(i >= setnamerange || j >= setnamerange || k >= setnamerange ||
        i < 0 || j < 0 || k < 0)
    {
        cout << "the set name is out of the range. " << endl;
        cout << "please use the set name among 0 to " << setnamerange - 1 << endl;
        return;
    }
    if(root[i] == EMPTY_ POSITION && root[j] == EMPTY_ POSITION)
    {
        cout << "error: the set i and j are both empty! " << endl;
        return;
    }
    int counti = (root[i] == EMPTY_ POSITION ? 0 : count[root[i]]);
    int countj = (root[j] == EMPTY_ POSITION ? 0 : count[root[j]]);

    // 为了保证 root[j]是大集合, root[i]是大集合
    if(counti > countj)
    {
        int t = i;
        i = j;
    }
}

```



```

        j = t;
    }
    int large = root[j];
    int small = root[i];
    if (small != EMPTY_POSITION) // 小的集合不是空集。
    {
        father[small] = large;
        name[small] = -1;
        count[large] = count[large] + count[small];
        count[small] = 0;
    }
    name[large] = k;
    root[i] = EMPTY_POSITION;
    root[j] = EMPTY_POSITION;
    root[k] = large;
}
void testUnionTreeSet()
{
    int root[] = {EMPTY_POSITION, 4, 7, 10, EMPTY_POSITION, EMPTY_POSITION, EMPTY_POSITION};
    int name[] = {EMPTY_POSITION, EMPTY_POSITION, EMPTY_POSITION, EMPTY_POSITION, 1,
                  EMPTY_POSITION, EMPTY_POSITION, 2, EMPTY_POSITION, EMPTY_POSITION, 3};
    int count[] = {0, 0, 0, 0, 5, 0, 0, 4, 0, 0, 2};
    int father[] = {EMPTY_POSITION, 6, 4, 2, EMPTY_POSITION, 4, 7,
                   EMPTY_POSITION, 6, 2, EMPTY_POSITION, 10};
    unionSet(father, name, count, root, 2, 3, 4, 7);
    cout << "新集合 4 中元素个数: " << count[root[4]] << endl;
    int setname = find(11, name, father);
    cout << "元素 11 在集合" << setname << "中" << endl;
    setname = find(8, name, father);
    cout << "元素 8 在集合" << setname << "中" << endl;
}
int main()
{
    testUnionTreeSet();
    system("pause");
    return 0;
}
int find(int v, int * name, int * father)
{
    list<int> l;
    l.clear();
    // 从顶点 v 向上追溯, 知道这个集合的根顶点。
    while (father[v] != -1)
    {
        l.push_back(v);
        v = father[v];
    }
    // v 此刻已经是根顶点。
    // 调整向上追溯的路径上的所有的顶点, 调整它们使得它们全部直接连在根上,
    // 减少之后搜索的代价。
    for (list<int>::iterator it = l.begin(); it != l.end(); it++)
    {
        father[*it] = v;
    }
    return name[v];
}

```

A. 18 求最小元素的离线算法

```
#include <iostream>
```

```

#include <list>
using namespace std;
int find(int v, int * name, int * father);
const int EMPTY_ POSITION = -1;
// fig 4.19 Executing instruction UNION(i, j, k)
// 集合求并算法实现(集合用树结构来表示)
// 输入:
//     father: 数组, 存放每个顶点的父顶点, 如 father[v]表示顶点 v 的父顶点
//     name: 数组, 存放每个根顶点对应的集合的名称
//     name[v]表示根顶点为 v 的集合的名称
//     count: 数组, 存放每个根顶点代表的集合中的元素的个数
//     count[v]表示根顶点是 v 的集合的元素的个数
//     root: 数组, 存放每个集合对应的根顶点。root[i]表示集合 i 的根顶点 v
//     i, j: 给外部名称为 i 和 j 的集合求并
//     k: 求并后的集合外部名称
//     setnamerange: 集合名的范围, 只能从 0 到 (setnamerange - 1)
// 输出:
//     运行结束后的数组中存放的就是求并之后的结果
// 测试函数:
//     void testUnionTreeSet()
// 测试数据:
//     求并之前的集合如下:
//     集合 1: 4      集合 2: 7      集合 3: 10
//           / \          \          /
//          2  5          6          11
//         / \          / \
//        3  9          1  8
// 将集合 2 和集合 3 求并, 求并后的新集合为 4。
// 测试结果:
//     集合 1 不变, 集合 4      7
//           / \
//          6  10
//         / \  \
//        1  8  11
void unionSet(int * father, int * name, int * count, int * root,
              int i, int j, int k, int setnamerange)
{
    if(i >= setnamerange || j >= setnamerange || k >= setnamerange ||
        i < 0 || j < 0 || k < 0)
    {
        cout << "the set name is out of the range. " << endl;
        cout << "please use the set name among 0 to " << setnamerange - 1 << endl;
        return;
    }
    if(root[i] == EMPTY_ POSITION && root[j] == EMPTY_ POSITION)
    {
        cout << "error: the set i and j are both empty! " << endl;
        return;
    }
    int counti = (root[i] == EMPTY_ POSITION ? 0 : count[root[i]]);
    int countj = (root[j] == EMPTY_ POSITION ? 0 : count[root[j]]);

    // 为了保证 root[j]是大集合, root[i]是大集合。
    if(counti > countj)
    {
        int t = i;
        i = j;
        j = t;
    }
    int large = root[j];

```

```

int small = root[i];
if (small != EMPTY_POSITION) // 小的集合不是空集
{
    father[small] = large;
    name[small] = -1;
    count[large] = count[large] + count[small];
    count[small] = 0;
}
name[large] = k;
root[i] = EMPTY_POSITION;
root[j] = EMPTY_POSITION;
root[k] = large;
}
void testUnionTreeSet()
{
    int root[] = {EMPTY_POSITION, 4, 7, 10, EMPTY_POSITION, EMPTY_POSITION, EMPTY_POSITION, 3, 4, 7};
    int name[] = {EMPTY_POSITION, EMPTY_POSITION, EMPTY_POSITION, EMPTY_POSITION, 1,
                  EMPTY_POSITION, EMPTY_POSITION, 2, EMPTY_POSITION, EMPTY_POSITION, 3};
    int count[] = {0, 0, 0, 0, 5, 0, 0, 4, 0, 0, 2};
    int father[] = {EMPTY_POSITION, 6, 4, 2, EMPTY_POSITION, 4, 7,
                    EMPTY_POSITION, 6, 2, EMPTY_POSITION, 10};
    unionSet(father, name, count, root, 2, 3, 4, 7);
    cout << "新集合 4 中元素个数: " << count[root[4]] << endl;
    int setname = find(11, name, father);
    cout << "元素 11 在集合" << setname << "中" << endl;
    setname = find(8, name, father);
    cout << "元素 8 在集合" << setname << "中" << endl;
}
int main()
{
    testUnionTreeSet();
    system("pause");
    return 0;
}
int find(int v, int * name, int * father)
{
    list<int> l;
    l.clear();
    // 从节点 v 向上追溯, 知道这个集合的根顶点
    while(father[v] != -1)
    {
        l.push_back(v);
        v = father[v];
    }
    // v 此刻已经是根顶点。
    // 调整向上追溯的路径上的所有的顶点, 调整它们使得它们全部直接连在根上,
    // 减少之后搜索的代价
    for(list<int>::iterator it = l.begin(); it != l.end(); it++)
    {
        father[*it] = v;
    }
    return name[v];
}

```

A. 19 2-3 树及相关算法

```

#include <iostream>
#include <queue>
using namespace std;

```

//2-3 树的顶点定义。

```

struct Node23Tree
{
    int lvalue; //左子树中最大的元素值
    int mvalue; //中子树中最大的元素值
    int maxValue; //以这个顶点为根的树中的最大的元素
    Node23Tree * sons[3];
    Node23Tree * father;
    Node23Tree(int leafValue)
    {
        this->lvalue = this->maxValue = this->mvalue = leafValue;
        father = NULL;
        for(int i=0; i<3; i++)
        {
            sons[i] = NULL;
        }
    }
    Node23Tree(int LValue, int MValue)
    {
        this->lvalue = LValue;
        this->mvalue = MValue;
        for(int i=0; i<3; i++)
        {
            sons[i] = NULL;
        }
        father = NULL;
    }
};

// 判断一个顶点是不是 2-3 树的叶顶点。
bool is23TreeLeaf(Node23Tree * node)
{
    if(node == NULL)
    {
        return false;
    }
    return (node->sons[0] == NULL &&
            node->sons[1] == NULL &&
            node->sons[2] == NULL &&
            node->lvalue == node->mvalue);
}

// 广度遍历打印出整棵树的所有顶点。
void show23Tree(Node23Tree * root)
{
    queue<Node23Tree * > q;
    q.push(root);
    while(! q.empty())
    {
        if(q.front() != NULL)
        {
            cout << q.front()->lvalue << ":" << q.front()->mvalue << " \t";
            q.push(q.front()->sons[0]);
            q.push(q.front()->sons[1]);
            q.push(q.front()->sons[2]);
        }
        q.pop();
    }
    cout << endl;
}

// 求一棵 2-3 树的高度
int Height23Tree(Node23Tree * root)

```

```

{
    int height = 0;
    Node23Tree * node = root;
    if (node == NULL)
    {
        return height;
    }
    while (true)
    {
        height ++;
        if (is23TreeLeaf (node)) // node 已经是叶顶点了
        {
            break; // 跳出 while 循环
        }
        if (node->sons[0] != NULL)
        {
            node = node->sons[0];
        }
        else if (node->sons[1] != NULL)
        {
            node = node->sons[1];
        }
        else
        {
            node = node->sons[2];
        }
    }
    return height;
}
// Fig 4.29. Procedure Search
// 2-3 树中的 search 函数
// 输入:
//     a: 被查找的值
//     root: 查找开始的根顶点
// 输出:
//     返回包含 a 值的叶顶点, 如果找不到符合条件的叶顶点,
//     返回这个值应该被插入的父顶点
// 测试函数:
//     void testSearch23tree();
// 测试数据:
//     在下列 2-3 树中查找值 5。非叶顶点(a: b), a 是这个顶点为根的左子树中最大的值,
//     b 是以这个顶点为根的中子树中最大的值
//           (3: 7)
//         /   |   \
//       (1: 3) (5: 6) (8: 9)
//     / | \ / | \ / | \
//    1 3 N 5 6 7 8 9 N
// 测试数据:
//     返回包含了 5 的叶顶点。
Node23Tree * search23tree(int a, Node23Tree * root)
{
    if (root == NULL)
    {
        return NULL;
    }
    for (int i = 0; i < 3; i++)
    {
        if (is23TreeLeaf (root->sons[i]) && root->sons[i]->lvalue == a)
        { //找到了含有 a 的叶顶点。
            return root->sons[i];
        }
    }
}

```

```

    }
}
if(is23TreeLeaf(root->sons[0]) ||
    is23TreeLeaf(root->sons[1]) ||
    is23TreeLeaf(root->sons[2]))
{
    return root; //在树中找不到 a
}
if(a <= root->lvalue)
{
    return search23tree(a, root->sons[0]);
}
else if(a <= root->mvalue)
{
    return search23tree(a, root->sons[1]);
}
else
{
    return search23tree(a, root->sons[2]);
}
}
void testSearch23tree()
{
    // 创建一棵 2-3 树以供测试
    Node23Tree * root = new Node23Tree(3, 7);
    root->sons[0] = new Node23Tree(1, 3);
    root->sons[1] = new Node23Tree(5, 6);
    root->sons[2] = new Node23Tree(8, 9);
    root->sons[0]->sons[0] = new Node23Tree(1, 1);
    root->sons[0]->sons[1] = new Node23Tree(3, 3);
    root->sons[1]->sons[0] = new Node23Tree(5, 5);
    root->sons[1]->sons[1] = new Node23Tree(6, 6);
    root->sons[1]->sons[2] = new Node23Tree(7, 7);
    root->sons[2]->sons[0] = new Node23Tree(8, 8);
    root->sons[2]->sons[1] = new Node23Tree(9, 9);
    Node23Tree * resultNode = search23tree(5, root);
    cout << resultNode->lvalue << endl;
    cout << resultNode->mvalue << endl;
    // 删除这棵 2-3 树所占用的空间
    delete root->sons[0]->sons[0];
    delete root->sons[0]->sons[1];
    delete root->sons[1]->sons[0];
    delete root->sons[1]->sons[1];
    delete root->sons[1]->sons[2];
    delete root->sons[2]->sons[0];
    delete root->sons[2]->sons[1];
    delete root->sons[0];
    delete root->sons[1];
    delete root->sons[2];
    delete root;
}
// fig 4.30.
//当 treeNode 的子结点已经满了之后, 向 treeNode 插入子结点时的算法
//输入:
// treeNode: 将待插入的值当成子结点插入到此结点。
// childToInsert: 待插入的结点(欲插入为 treeNode 的子结点,
// 但是会造成 treeNode 的子结点超过 3, 需要调整)
// insertPosition: 插入到 treeNode 中子结点的位置,
// 分别是: 0, 1, 2, 3;
//输出:

```

```

// 调整插入完成之后的树， 返回的是插入后的新树的根， 如果没有产生新的根，
// 则返回 NULL
//测试函数：
// void testAddSon()
//测试数据：
// 在下列 2-3 树中 (5: 6) 结点中插入 4。 非叶结点 (a: b)，
// a 是这个结点为根的左子树中最大的值，
// b 是以这个结点为根的中子树中最大的值
//          (3: 7)
//         /  |  \
//       (1: 3) (5: 6) (8: 9)
//      / | \ / | \ / | \
//     1 3 N 5 6 7 8 9 N
// 测试数据：
// 在 (5: 6) 结点中插入 4 之后的 2-3 树如下：
//          (5: 9)
//         /  \
//       (3: 5) (7: 9)
//      /  \  /  \
//    (1: 3) (4: 5) (6: 7) (8: 9)
//   /  \ /  \ /  \ /  \
//  1 3 4 5 6 7 8 9
Node23Tree * addSon(Node23Tree * treeNode, Node23Tree * childToInsert, int insertPosition)
{
    Node23Tree * pNodes [4];
    for(int i=0; i<4; i++)
    {
        if(i<insertPosition)
        {
            pNodes[i] = treeNode->sons[i];
        }
        else if(i == insertPosition)
        {
            pNodes[i] = childToInsert;
        }
        else
        {
            pNodes[i] = treeNode->sons[i-1];
        }
    }
    Node23Tree * newNode = new Node23Tree(pNodes[2] ->maxValue, pNodes[3] ->maxValue);
    newNode->maxValue = pNodes[3] ->maxValue;
    treeNode->sons[0] = pNodes[0];
    treeNode->lvalue = pNodes[0] ->maxValue;
    pNodes[0] ->father = treeNode;
    treeNode->sons[1] = pNodes[1];
    treeNode->mvalue = pNodes[1] ->maxValue;
    pNodes[1] ->father = treeNode;
    treeNode->sons[2] = NULL;
    treeNode->maxValue = pNodes[1] ->maxValue;
    newNode->sons[0] = pNodes[2];
    pNodes[2] ->father = newNode;
    newNode->sons[1] = pNodes[3];
    pNodes[3] ->father = newNode;
    if(treeNode->father == NULL) // treeNode 是根结点， 没有父结点。
    {
        Node23Tree * root = new Node23Tree(treeNode->maxValue, newNode->maxValue);
        root->sons[0] = treeNode;
        treeNode->father = root;
        root->sons[1] = newNode;
    }
}

```

```

    treeNode->father = root;
    return root;
}
else // treeNode 有父节点。
{
    if (treeNode->father->sons[0] == NULL ||
        treeNode->father->sons[1] == NULL ||
        treeNode->father->sons[2] == NULL)
    {
        // father 的子结点没有满，还可以直接插入
        if (treeNode == treeNode->father->sons[0])
        {
            if (treeNode->father->sons[1] == NULL)
            {
                newNode->father = treeNode->father;
                treeNode->father->sons[1] = newNode;
                treeNode->father->mvalue = newNode->maxValue;
            }
            else
            {
                treeNode->father->sons[2] = treeNode->father->sons[1];
                treeNode->father->sons[1] = newNode;
                treeNode->father->mvalue = newNode->maxValue;
                newNode->father = treeNode->father;
            }
        }
        else if (treeNode->father->sons[1] == treeNode)
        {
            if (treeNode->father->sons[2] == NULL)
            {
                treeNode->father->sons[2] = newNode;
                newNode->father = treeNode->father;
            }
            else
            {
                treeNode->father->sons[0] = treeNode->father->sons[1];
                treeNode->father->sons[1] = treeNode->father->sons[2];
                treeNode->father->sons[2] = newNode;
                newNode->father = treeNode->father;
            }
            treeNode->father->mvalue = newNode->maxValue;
        }
        else // (treeNode == treeNode->father->sons[2])
        {
            if (treeNode->father->sons[0] == NULL)
            {
                treeNode->father->sons[0] = treeNode->father->sons[1];
                treeNode->father->sons[1] = treeNode->father->sons[2];
                treeNode->father->sons[2] = newNode;
                newNode->father = treeNode->father;
            }
            else // treeNode->father->sons[1] == NULL
            {
                treeNode->father->sons[1] = treeNode->father->sons[2];
                treeNode->father->sons[2] = newNode;
                newNode->father = treeNode->father;
            }
            treeNode->father->mvalue = newNode->maxValue;
        }
    }
    return NULL;
}

```



```

    }
    else // treeNode -> father 的子结点完全已经满了，必须递归调用 addSon 调整插入
    {
        // insertPos 是插入 newNode 的位置，应该刚好是 treeNode 的右一个结点。
        int insertPos = 0;
        if (treeNode == treeNode -> father -> sons[0])
        {
            insertPos = 1;
        }
        else if (treeNode == treeNode -> father -> sons[1])
        {
            insertPos = 2;
        }
        else // treeNode == treeNode -> father -> sons[2]
        {
            insertPos = 3;
        }
        return addSon(treeNode -> father, newNode, insertPos);
    }
}

void testAddSon()
{
    // 创建一棵 2-3 树以供测试
    Node23Tree * root = new Node23Tree(3, 7);
    root -> maxValue = 9;
    root -> sons[0] = new Node23Tree(1, 3);
    root -> sons[0] -> maxValue = 3;
    root -> sons[0] -> father = root;
    root -> sons[1] = new Node23Tree(5, 6);
    root -> sons[1] -> maxValue = 7;
    root -> sons[1] -> father = root;
    root -> sons[2] = new Node23Tree(8, 9);
    root -> sons[2] -> maxValue = 9;
    root -> sons[2] -> father = root;
    root -> sons[0] -> sons[0] = new Node23Tree(1);
    root -> sons[0] -> sons[0] -> father = root -> sons[0];
    root -> sons[0] -> sons[1] = new Node23Tree(3);
    root -> sons[0] -> sons[1] -> father = root -> sons[0];
    root -> sons[1] -> sons[0] = new Node23Tree(5);
    root -> sons[1] -> sons[0] -> father = root -> sons[1];
    root -> sons[1] -> sons[1] = new Node23Tree(6);
    root -> sons[1] -> sons[1] -> father = root -> sons[1];
    root -> sons[1] -> sons[2] = new Node23Tree(7);
    root -> sons[1] -> sons[2] -> father = root -> sons[1];
    root -> sons[2] -> sons[0] = new Node23Tree(8);
    root -> sons[2] -> sons[0] -> father = root -> sons[2];
    root -> sons[2] -> sons[1] = new Node23Tree(9);
    root -> sons[2] -> sons[1] -> father = root -> sons[2];
    show23Tree(root);
    cout << "the height of the tree is : " << Height23Tree(root) << endl;
    Node23Tree * curNode = root -> sons[1];
    Node23Tree * newNode = new Node23Tree(4);
    Node23Tree * newRoot = addSon(curNode, newNode, 0);
    if (newRoot != NULL)
    {
        root = newRoot;
    }

    show23Tree(root);
}

```

```

    cout << "the height of the tree is : " << Height23Tree(root) << endl;
}
//Fig. 4.32 Procedure IMPLANT
// 功能: 合并两棵 2-3 树
// 输入:
//      tree1: 待合并的第一棵树的根结点
//      tree2: 待合并的第二棵树的根结点
// (两棵树必须满足条件都是 2-3 树且 tree2 的所有结点值都比 tree1 中大)
// 输出:
//      返回合并后的新的树的根结点
// 测试函数:
//      void testImplant23Tree()
// 测试数据:      tree1:                      tree2:
//                (3: 7)                      (10: 12)
//                / | \                      / | \
//                (1: 3) (5: 6) (8: 9)        10 12 13
//                / | \ / | \ / | \
//                1 3 N 5 6 7 8 9 N
// 测试结果:
//                (7: 13)
//                /      \
//                (3: 7)   (9: 13)
//                / \     / \
//                (1: 3) (5: 6) (8: 9) (10: 12)
//                / \ / | \ / \ / \
//                1 3 5 6 7 8 9 10 12 13
Node23Tree * implant23Tree(Node23Tree * tree1, Node23Tree * tree2)
{
    int h1 = Height23Tree(tree1);
    int h2 = Height23Tree(tree2);
    if (h1 == h2) // 两棵树的高度一样
        //使其同时成为一个新的根结点的子结点之后仍然可以保持平衡
    {
        Node23Tree * root = new Node23Tree(tree1->maxValue, tree2->maxValue);
        return root; // 完成合并
    }
    else
    {
        // 保证 t1 的高度比 t2 要大
        Node23Tree * t1 = NULL;
        Node23Tree * t2 = NULL;
        bool isreversed = false;
        if (h1 > h2)
        {
            t1 = tree1;
            t2 = tree2;
        }
        else
        {
            isreversed = true;
            t1 = tree2;
            t2 = tree1;
            int temp = h1;
            h1 = h2;
            h2 = temp;
        }
        int height = 1;
        Node23Tree * node = t1;
        int rightMost = 2;
        while(true)

```

```

{
    if (height >= h1 - h2)
    {
        break;
    }
    // 如果 tree1 和 tree2 位置没有反过来, 则沿着 t1 最右边的路径, 向下一层
    // 如果 tree1 和 tree2 位置反了过来, 则沿着 t1 最左边的路径, 向下一层
    if (isreversed)
    {
        rightMost = 0;
    }
    if (t1 -> sons[rightMost] != NULL)
    {
        node = t1 -> sons[rightMost];
    }
    else
    {
        node = t1 -> sons[1];
    }
    height++;
}
// node 现在就是 t2 要插入的位置, 将 t2 的根结点作为 node 的最右的子结点插入
// (如果 tree1 和 tree2 曾经调换位置, 则是最左的子结点)
if (node -> sons[0] != NULL &&
    node -> sons[1] != NULL &&
    node -> sons[2] != NULL)
{
    // node 所有的子结点已经满了, 需要调整插入
    if (isreversed)
    {
        return addSon(node, t2, 0);
    }
    else
    {
        return addSon(node, t2, 3);
    }
}
// node 还有空的子结点可以插入
int leftMost = 0;
if (isreversed)
{
    leftMost = 2;
}
if (node -> sons[leftMost] == NULL)
{
    node -> sons[leftMost] = node -> sons[1];
    node -> sons[1] = node -> sons[rightMost];
}
else if (node -> sons[1] == NULL)
{
    node -> sons[1] = node -> sons[rightMost];
    node -> mvalue = node -> sons[rightMost] -> maxValue;
}
// 此刻 node -> sons[rightmost] 一定为空, 可以插入
node -> sons[rightMost] = t2;
if (isreversed) // leftmost = 2; rightmost = 0;
{
    node -> mvalue = node -> sons[1] -> maxValue;
    node -> lvalue = t2 -> maxValue;
    node -> maxValue = node -> sons[2] -> maxValue;
}

```

```

    }
    else // rightmost = 2; leftmost = 0;
    {
        node -> maxValue = t2 -> maxValue;
        node -> lvalue = node -> sons[0] -> maxValue;
        node -> mvalue = node -> sons[1] -> maxValue;
    }
    t2 -> father = node;
    // 向上回溯一直到根, 更新 maxValue 值
    while (node != NULL)
    {
        int max = node -> maxValue;
        node = node -> father;
        if (node != NULL)
        {
            node -> maxValue = max;
        }
    }
    return t1;
}

void testImplant23Tree()
{
    // 创建 2-3 树 tree1 以供测试
    Node23Tree * tree1 = new Node23Tree(3, 7);
    tree1 -> maxValue = 9;
    tree1 -> sons[0] = new Node23Tree(1, 3);
    tree1 -> sons[0] -> maxValue = 3;
    tree1 -> sons[0] -> father = tree1;
    tree1 -> sons[1] = new Node23Tree(5, 6);
    tree1 -> sons[1] -> maxValue = 7;
    tree1 -> sons[1] -> father = tree1;
    tree1 -> sons[2] = new Node23Tree(8, 9);
    tree1 -> sons[2] -> maxValue = 9;
    tree1 -> sons[2] -> father = tree1;
    tree1 -> sons[0] -> sons[0] = new Node23Tree(1);
    tree1 -> sons[0] -> sons[0] -> father = tree1 -> sons[0];
    tree1 -> sons[0] -> sons[1] = new Node23Tree(3);
    tree1 -> sons[0] -> sons[1] -> father = tree1 -> sons[0];
    tree1 -> sons[1] -> sons[0] = new Node23Tree(5);
    tree1 -> sons[1] -> sons[0] -> father = tree1 -> sons[1];
    tree1 -> sons[1] -> sons[1] = new Node23Tree(6);
    tree1 -> sons[1] -> sons[1] -> father = tree1 -> sons[1];
    tree1 -> sons[1] -> sons[2] = new Node23Tree(7);
    tree1 -> sons[1] -> sons[2] -> father = tree1 -> sons[1];
    tree1 -> sons[2] -> sons[0] = new Node23Tree(8);
    tree1 -> sons[2] -> sons[0] -> father = tree1 -> sons[2];
    tree1 -> sons[2] -> sons[1] = new Node23Tree(9);
    tree1 -> sons[2] -> sons[1] -> father = tree1 -> sons[2];
    show23Tree(tree1);
    // 创建 2-3 树 tree2 以供测试
    Node23Tree * tree2 = new Node23Tree(10, 12);
    tree2 -> maxValue = 13;
    tree2 -> sons[0] = new Node23Tree(10);
    tree2 -> sons[1] = new Node23Tree(12);
    tree2 -> sons[2] = new Node23Tree(13);
    for (int i = 0; i < 3; i++)
    {
        tree2 -> sons[i] -> father = tree2;
    }
}

```

```

    show23Tree(tree2);
    Node23Tree * resultRoot = implant23Tree(tree1, tree2);
    show23Tree(resultRoot);
}
// Fig. 4.34 Procedure to split a 2-3 tree.
// 2-3 树的分裂算法。
// 输入:
//     value: 分裂 2-3 树的临界值, 将 < - value 的结点放在分裂后的 一棵子树中,
//     将 > value 的结点放在分裂后的另一棵子树中
//     tree: 被分裂的 2-3 树的根结点
// 输出:
//     subtree1: 分裂成两棵树后的第一棵树(包括了 < = value 的结点)的根结点
//     subtree2: 分裂成两棵树后的第二棵树(包括了 > value 的结点)的根结点
// 测试函数:
//     void testDevide23Tree()
// 测试数据:
//     被分裂的 2-3 树如下所示: 分裂的临界值是 6
//           (7: 13)
//         /      \
//       (3: 7)    (9: 13)
//      /  \      /  \
//    (1: 3) (5: 6) (8: 9) (10: 12)
//   /  \  /  \  /  \  /  \
//  1   3 5   6 7   8   9  10 12 13
// 运行结果:
//     分裂后的两棵树分别如下。
//           (3: 6)                (9: 13)
//         /      \              /      \
//       (1: 3)    (5: 6)        (7: 8)   (10: 12)
//      /  \      /  \        /  |  \  /  |  \
//     1   3   5   6        7   8   9  10 12 13
void devide23Tree(int value, Node23Tree * tree, Node23Tree * & subtree1, Node23Tree * & subtree2)
{
    queue<Node23Tree * > treelessequalvalue;
    // 存放删除后所有结点小于等于 value 的树的根
    queue<Node23Tree * > treegreatervalue; // 存放删除后所有结点大于等于 value 的树的根
    Node23Tree * node = tree;
    // 从 tree 的根结点一直到值为 value 的叶结点, 将这条路径上的所有的结点都删除
    while(! is23TreeLeaf(node) || node == NULL)
    {
        Node23Tree * back = node;
        int sonIdx = 0;
        // 移动到路径的下一个结点。
        if(value <= node->lvalue)
        {
            sonIdx = 0;
        }
        else if(value <= node->mvalue)
        {
            sonIdx = 1;
        }
        else
        {
            sonIdx = 2;
        }
        node = node->sons[sonIdx];
        for(int i = 0; i < 3; i++)
        {
            if(back->sons[i] != NULL)

```

```

        {
            back -> sons[i] -> father = NULL;
        }
    }
    for(int i=0; i < sonIdx; i++)
    {
        if(back -> sons[i] != NULL)
        {
            treelessequalvalue.push(back -> sons[i]);
        }
    }
    if(is23TreeLeaf(node)) // 是叶结点
    {
        treelessequalvalue.push(node);
        // 叶结点不会被删除, 所以作为一个子树放到临时存储空间
    }
    for(int i = sonIdx+1; i < 3; i++)
    {
        if(back -> sons[i] != NULL)
        {
            treegreatervalue.push(back -> sons[i]);
        }
    }
    delete back; // 删除这条路径中的当前结点(不是叶结点)
}
while(treegreatervalue.size() > 1)
{
    // 将这些子树合并
    Node23Tree * node1 = treegreatervalue.front();
    treegreatervalue.pop();
    Node23Tree * node2 = treegreatervalue.front();
    treegreatervalue.pop();
    if(node1 -> lvalue <= node2 -> lvalue)
    {
        treegreatervalue.push(implant23Tree(node1, node2));
    }
    else
    {
        treegreatervalue.push(implant23Tree(node2, node1));
    }
}
subtree2 = treegreatervalue.front();
treegreatervalue.pop();
while(treelessequalvalue.size() > 1)
{
    // 将这些子树合并
    Node23Tree * node1 = treelessequalvalue.front();
    treelessequalvalue.pop();
    Node23Tree * node2 = treelessequalvalue.front();
    treelessequalvalue.pop();
    if(node1 -> lvalue <= node2 -> lvalue)
    {
        treelessequalvalue.push(implant23Tree(node1, node2));
    }
    else
    {
        treelessequalvalue.push(implant23Tree(node2, node1));
    }
}
subtree1 = treelessequalvalue.front();

```

```

    treelessequalvalue.pop();
}
void testDevide23Tree()
{
    // 建立一棵 2-3 树以供测试
    Node23Tree * root = new Node23Tree(7, 13);
    root->maxValue = 13;
    root->sons[0] = new Node23Tree(3, 7);
    root->sons[0]->maxValue = 7;
    root->sons[0]->father = root;
    root->sons[1] = new Node23Tree(9, 13);
    root->sons[1]->maxValue = 13;
    root->sons[1]->father = root;
    root->sons[0]->sons[0] = new Node23Tree(1, 3);
    root->sons[0]->sons[0]->maxValue = 3;
    root->sons[0]->sons[0]->father = root->sons[0];
    root->sons[0]->sons[1] = new Node23Tree(5, 6);
    root->sons[0]->sons[1]->maxValue = 7;
    root->sons[0]->sons[1]->father = root->sons[0];
    root->sons[1]->sons[0] = new Node23Tree(8, 9);
    root->sons[1]->sons[0]->maxValue = 9;
    root->sons[1]->sons[0]->father = root->sons[1];
    root->sons[1]->sons[1] = new Node23Tree(10, 12);
    root->sons[1]->sons[1]->maxValue = 13;
    root->sons[1]->sons[1]->father = root->sons[1];
    root->sons[0]->sons[0]->sons[0] = new Node23Tree(1);
    root->sons[0]->sons[0]->sons[0]->father = root->sons[0]->sons[0];
    root->sons[0]->sons[0]->sons[1] = new Node23Tree(3);
    root->sons[0]->sons[0]->sons[1]->father = root->sons[0]->sons[0];
    root->sons[0]->sons[1]->sons[0] = new Node23Tree(5);
    root->sons[0]->sons[1]->sons[0]->father = root->sons[0]->sons[1];
    root->sons[0]->sons[1]->sons[1] = new Node23Tree(6);
    root->sons[0]->sons[1]->sons[1]->father = root->sons[0]->sons[1];
    root->sons[0]->sons[1]->sons[2] = new Node23Tree(7);
    root->sons[0]->sons[1]->sons[2]->father = root->sons[0]->sons[1];
    root->sons[1]->sons[0]->sons[0] = new Node23Tree(8);
    root->sons[1]->sons[0]->sons[0]->father = root->sons[1]->sons[0];
    root->sons[1]->sons[0]->sons[1] = new Node23Tree(9);
    root->sons[1]->sons[0]->sons[1]->father = root->sons[1]->sons[0];
    root->sons[1]->sons[1]->sons[0] = new Node23Tree(10);
    root->sons[1]->sons[1]->sons[0]->father = root->sons[1]->sons[1];
    root->sons[1]->sons[1]->sons[1] = new Node23Tree(12);
    root->sons[1]->sons[1]->sons[1]->father = root->sons[1]->sons[1];
    root->sons[1]->sons[1]->sons[2] = new Node23Tree(13);
    root->sons[1]->sons[1]->sons[2]->father = root->sons[1]->sons[1];
    show23Tree(root);
    Node23Tree * subroot1 = NULL;
    Node23Tree * subroot2 = NULL;
    devide23Tree(6, root, subroot1, subroot2);
    cout << "come here! " << endl;
    show23Tree(subroot1);
    show23Tree(subroot2);
}
int main()
{
    testSearch23tree();
    testAddSon();
    testImplant23Tree();
    testDevide23Tree();
    system("pause");
}

```

```
    return 0;
}
```

A.20 集合划分的算法

```
#include <iostream>
#include <vector>
#include <set>
using namespace std;

int f_ inverse(int input, int setsize)
{
    return (input + 1) % setsize;
}

//Fig 4.35 Partitioning algorithm
// 集合划分的算法实现
// 输入:
//    blocks: 所有集合的列表
// 输出:
//    划分后的集合
// 测试函数:
//    void testPartitionSet()
//    根据 f_ inverse 规则来划分集合
// 测试数据:
//    block1 = {0, 1, 2, 3, 4, 5}; block2 = { 6, 7 }; block3 = {8};
// 测试结果:
//    block1 = {4, 5}; block2 = {6, 7}; block3 = {8}; block4 = {0}; block5 = {1}; block6 = {2, 3};
void partitionSet(vector<set<int>> &blocks)
{
    vector<int> waitinglist;
    for(int i=0; i<blocks.size(); i++)
    {
        waitinglist.push_back(i);
    }
    int endBlockID=blocks.size()-1;
    set<int> inverse;
    while(! waitinglist.empty())
    {
        int i=waitinglist.back();
        waitinglist.pop_back();
        for(set<int>:: iterator it=blocks.at(i).begin();
            it != blocks.at(i).end();
            it++)
        {
            inverse.insert(f_ inverse(* it, blocks.at(i).size()));
        }
        // 对于 blocks 里面的每一个 block
        for(int j=0; j<=endBlockID; j++)
        {
            set<int> newblock;
            newblock.clear();
            for(set<int>:: iterator iit=blocks.at(j).begin();
                iit != blocks.at(j).end();
                )
            {
                if(inverse.find(* iit) != inverse.end())
                {
                    // blocks[j] 交 inverse 不为空
                    newblock.insert(* iit); // 新的集合里面放 block[j]交 inverse 的元素
                    set<int>:: iterator deleteit=iit;
                }
            }
        }
    }
}
```



```

        iit++;
        blocks.at(j).erase(deleteit);
    }
    else
    {
        iit++;
    }
}
if(blocks.at(j).size() > 0 // blocks[j] 不是 inverse 的子集
    && newblock.size() > 0) // blocks[j] 与 inverse 的交集不是空集
{
    blocks.push_back(newblock);
    endBlockID++;
    vector<int>::iterator vecIt = waitinglist.begin();
    for(vecIt = waitinglist.begin(); vecIt != waitinglist.end(); vecIt++)
    {
        if(*vecIt == j) // j 在等待队列中
        {
            waitinglist.push_back(endBlockID);
            break;
        }
    }
    if(vecIt == waitinglist.end()) // j 不在等待队列中;
    {
        if(blocks.at(j).size() <= newblock.size())
        {
            waitinglist.push_back(j);
        }
        else
        {
            waitinglist.push_back(endBlockID);
        }
    }
}
if(blocks.at(j).size() == 0)
{
    blocks[j] = newblock;
}
}
}

void testPartitionSet()
{
    set<int> block1;
    for(int i=0; i<6; i++)
    {
        block1.insert(i);
    }
    set<int> block2;
    block2.insert(6);
    block2.insert(7);
    set<int> block3;
    block3.insert(8);
    vector<set<int>> blocks;
    blocks.push_back(block1);
    blocks.push_back(block2);
    blocks.push_back(block3);

    int i=0;
    for(vector<set<int>>::iterator it=blocks.begin(); it != blocks.end(); it++)

```

```

{
    cout << "set" << i++ << ": " << endl;
    for(set <int >:: iterator setit = it ->begin(); setit != it ->end(); setit++)
    {
        cout << * setit << " \t";
    }
    cout << endl;
}
partitionSet(blocks);
i=0;
for(vector <set <int > >:: iterator it = blocks.begin(); it != blocks.end(); it++)
{
    cout << "set" << i++ << ": " << endl;
    for(set <int >:: iterator setit = it ->begin(); setit != it ->end(); setit++)
    {
        cout << * setit << " \t";
    }
    cout << endl;
}
}
int main()
{
    testPartitionSet();
    system("pause");
    return 0;
}

```

A. 21 最小生成树算法

```

#include <iostream>
#include <list>
#include <vector>
#include <algorithm>
#include <stack>
#include <stdlib.h>
using namespace std;
struct Edge
{
    int start;
    int end;
    float weight;
    Edge(int _ start, int _ end, float _ weight)
    {
        start = _ start;
        end = _ end;
        weight = _ weight;
    }
    bool operator < (const Edge & e) const
    {
        return this ->weight < e.weight;
    }
    bool operator > (const Edge & e) const
    {
        return this ->weight > e.weight;
    }
    bool operator >= (const Edge & e) const
    {
        return this ->weight >= e.weight;
    }
    bool operator <= (const Edge & e) const

```

```

    {
        return this->weight <= e.weight;
    }
    bool operator == (const Edge & e) const
    {
        return (this->weight == e.weight &&
            this->start == e.start &&
            this->end == e.end);
    }
};

int find(int v, int * name, int * father);
void unionSet(int * father, int * name, int * count, int * root,
    int i, int j, int k, int setnamerange);

// Fig. 5.2
// 最小生成树的生成
// 输入:
//     V: 原图的点集(点从1到nV-1)
//     nV: 原图的点的数目
//     E: 原图的边集
// 输出:
//     生成的最小生成树边集为T, 点集仍然为V
// 测试函数:
//     testBuildMinimunCostSpanningTree()
// 测试数据:
//     原图由4点组成: 0, 1, 2, 3, 4, 5, 6
//
//          20
//          v0 - - - - - v1
//          / \       / \
//          23/  \1   4/  \15
//          / 36  \  / 9   \
//          v5 - - - - - v6 - - - - - v2
//          \   /   \   /
//          28 \ /25  \16 /3
//          \ / 17  \ /
//          v4 - - - - - v3
// 测试结果:
//     生成树点集与原图一致: 0, 1, 2, 3, 4, 5, 6
//
//          v0          v1
//          / \       /
//          23/  \1   4/
//          /   \  / 9
//          v5          v6 - - - - - v2
//                      /
//                      /3
//                  17 /
//          v4 - - - - - v3
void buildMinimunCostSpanningTree(int * V, int nV, vector <Edge> & E, vector <Edge> & T)
{
    T.clear();
    // 初始化一个集合的集合
    int * father = new int[nV];
    int * root = new int[nV];
    int * count = new int[nV];
    int * name = new int[nV];
    // 每个点都被初始化为一个单独的点集。
    for(int i=0; i<nV; i++)
    {
        father[i] = -1;
        count[i] = 1; // 每个集合i中都只有i这一个元素
    }
}

```

```

    root[i] = i; // 集合 i 的根节点就是点 i;
    name[i] = i;
}
// 用堆的形式来存放图中所有的边 E;
make_heap(E.begin(), E.end(), greater<Edge>());
vector<Edge>::iterator edgeEndIt = E.end();
while(true)
{
    if(count[root[find(0, name, father)]] >= nV || edgeEndIt == E.begin())
    { // 所有的点都在同一个点集中了, 表示 T 中已经覆盖了所有的点
        break;
    }
    // 从 E 中选出权重最小的一条边 (v, w);
    Edge e = E.at(0);
    // 从 E 中删除 (v, w);
    pop_heap(E.begin(), edgeEndIt, greater<Edge>());
    edgeEndIt--;
    int w1 = find(e.start, name, father);
    int w2 = find(e.end, name, father);
    if(w1 != w2) // v 和 w 在 VS 中的不同点集 w1 和 w2 中
    {
        // 用点集 w1 union w2 代替 VS 中的点集 w1 和 w2;
        unionSet(father, name, count, root, w1, w2, w2, nV);
        // 将边 (v, w) 加入到 T 中;
        T.push_back(e);
    }
}
}

void testBuildMinimumCostSpanningTree()
{
    int V[] = {0, 1, 2, 3, 4, 5, 6};
    vector<Edge> E;
    Edge e01(0, 1, 20);
    Edge e12(1, 2, 15);
    Edge e23(2, 3, 3);
    Edge e34(3, 4, 17);
    Edge e45(4, 5, 28);
    Edge e50(5, 0, 23);
    Edge e06(0, 6, 1);
    Edge e16(1, 6, 4);
    Edge e26(2, 6, 9);
    Edge e36(3, 6, 16);
    Edge e46(4, 6, 25);
    Edge e56(5, 6, 36);
    E.push_back(e01);
    E.push_back(e12);
    E.push_back(e23);
    E.push_back(e34);
    E.push_back(e45);
    E.push_back(e50);
    E.push_back(e06);
    E.push_back(e16);
    E.push_back(e26);
    E.push_back(e36);
    E.push_back(e46);
    E.push_back(e56);
    vector<Edge> tree;
    tree.clear();
    buildMinimumCostSpanningTree(V, 7, E, tree);
    cout << "In the min cost spanning tree built: " << endl;
}

```

```

    for(vector<Edge>:: iterator it = tree.begin(); it != tree.end(); it++)
    {
        cout << "(" << it->start << ", " << it->end << "): " << it->weight << endl;
    }
}

int main()
{
    testBuildMinimumCostSpanningTree();
    system("pause");
    return 0;
}

// Fig 4.18 Executing instruction FIND(i)
// 在集合中查找某个元素所在的集合名(集合用树的结构来表示)
// 输入:
//     v : 查找结点 v
//     name: 数组, 存放每个根结点对应的集合的名称
//     name[v]表示根结点为 v 的集合的名称
//     father: 数组, 存放每个结点的父结点, 如 father[v]表示结点 v 的父结点
// 输出:
//     Vi 所属于的集合名将被打印到控制台(console)上, 并返回 Vi 所属于的集合名
// 测试数据: 给定集合:
//           7
//          / \
//         5  2
//        /
//       1
//      /
//     3
// 查找 v = 3;
// 测试结果: 找出根结点 7, 打印出集合名 3, 并将集合调整为一下结构
//           7
//          / \ \ \
//         5  2 1 3
int find(int v, int * name, int * father)
{
    list<int> l;
    l.clear();
    // 从节点 v 向上追溯, 知道这个集合的根结点。
    while(father[v] != -1)
    {
        l.push_back(v);
        v = father[v];
    }
    // v 此刻已经是根结点
    // 调整向上追溯的路径上的所有的结点, 调整它们使得它们全部直接连在根上,
    // 减少之后搜索的代价
    for(list<int>:: iterator it = l.begin(); it != l.end(); it++)
    {
        father[* it] = v;
    }
    return name[v];
}

// fig 4.19 Executing instruction UNION(i, j, k)
// 集合求并算法实现(集合用树结构来表示)
// 输入:
//     father: 数组, 存放每个结点的父结点, 如 father[v]表示结点 v 的父结点
//     name: 数组, 存放每个根结点对应的集合的名称
//     name[v]表示根结点为 v 的集合的名称
//     count: 数组, 存放每个根结点代表的集合中的元素的个数
//     count[v]表示根结点是 v 的集合的元素的个数

```

```

//    root: 数组, 存放每个集合对应的根结点。root[i]表示集合 i 的根结点 v
//    i, j: 给外部名称为 i 和 j 的集合求并
//    k: 求并后的集合外部名称
//    setnamerange: 集合名的范围, 只能从 0 到 (setnamerange - 1)
// 输出:
//    运行结束后的数组中存放的就是求并之后的结果
// 测试数据:
//    求并之前的集合如下:
//    集合 1: 4      集合 2: 7      集合 3: 10
//           / \      \      /
//          2  5      6      11
//         / \      / \
//        3  9      1  8
//    将集合 2 和集合 3 求并, 求并后的新集合为 4
// 测试结果:
//    集合 1 不变, 集合 4      7
//           / \
//          6  10
//         / \ \
//        1  8  11
void unionSet(int * father, int * name, int * count, int * root, int i, int j, int k, int setnamerange)
{
    if(i >= setnamerange || j >= setnamerange || k >= setnamerange ||
       i < 0 || j < 0 || k < 0)
    {
        cout << "the set name is out of the range. " << endl;
        cout << "please use the set name among 0 to " << setnamerange - 1 << endl;
        return;
    }
    if(root[i] == -1 && root[j] == -1)
    {
        cout << "error: the set i and j are both empty! " << endl;
        return;
    }
    int counti = root[i] == -1 ? 0 : count[root[i]];
    int countj = root[j] == -1 ? 0 : count[root[j]];
    if(counti > countj)
    {
        int t = i;
        i = j;
        j = t;
    }
    int large = root[j];
    int small = root[i];
    if(small != -1) // 小的集合不是空集
    {
        father[small] = large;
        name[small] = -1;
        count[small] = 0;
    }
    count[large] = count[large] + count[small];
    name[large] = k;
    root[i] = -1;
    root[j] = -1;
    root[k] = large;
}

```

A.22 图的深度优先遍历算法

```
#include <iostream>
```

```

#include <vector>
using namespace std;
struct Edge
{
    int start;
    int end;
    float weight;
    Edge(int _start, int _end, float _weight)
    {
        start = _start;
        end = _end;
        weight = _weight;
    }
    bool operator < (const Edge & e) const
    {
        return this->weight < e.weight;
    }
    bool operator > (const Edge & e) const
    {
        return this->weight > e.weight;
    }
    bool operator >= (const Edge & e) const
    {
        return this->weight >= e.weight;
    }
    bool operator <= (const Edge & e) const
    {
        return this->weight <= e.weight;
    }
    bool operator == (const Edge & e) const
    {
        return (this->weight == e.weight &&
            this->start == e.start &&
            this->end == e.end);
    }
};

// fig 5.6
// 图的深度优先遍历的算法实现
// 输入:
//     vertex: 图中开始搜索的顶点
//     adjacentList: 数组, 存放所有顶点的邻接点列表
//     adjacentList[i]表示顶点 i 的所有邻接点的列表
//     isOld: 布尔值数组, 存放一个顶点是否已经被访问过了,
//           如果顶点[i]已经访问过则 isOld[i]值为 true, 否则值为 false;
//     nVertex: 图中顶点的数目, 顶点编号从 0 到 (nVertex-1);
// 输出:
//     tree: 深度优先树的边集 (点集和原来的图的点集一致)
// 测试函数:
//     void testSearchDepthFirst()
// 测试数据:
//     遍历的起点是 v0
//     原图如下:
//           (v3)
//           | \
//           |  \
//           |   \
//     (v4) - - (v0) - - (v1)
//           |   /  \   |
//           |   /   \   |

```

```

//          | /      \ |
//          (v5)      (v2)
// 测试结果:
//      生成的深度优先遍历树: (树的根是 v0)
//          (v3)
//          \
//          \
//          \
//      (v4) -- (v0) -- (v1)
//          |           |
//          |           |
//          |           |
//      (v5)      (v2)

void searchDepthFirst(int vertex, vector<int> * adjacentList, bool * isOld, int nVertex,
vector<Edge> & tree)
{
    if(vertex >= nVertex)
    {
        return;
    }
    isOld[vertex] = true;

    for(vector<int>::iterator it = adjacentList[vertex].begin(); it != adjacentList[vertex].end(); it++)
    {
        if(isOld[* it] == false)
        {
            tree.push_back(Edge(vertex, * it, 1));
            searchDepthFirst(* it, adjacentList, isOld, nVertex, tree);
        }
    }
}

void testSearchDepthFirst()
{
    vector<Edge> dfTree;
    vector<int> adjacentList[6];
    adjacentList[0].push_back(0);
    adjacentList[0].push_back(1);
    adjacentList[0].push_back(2);
    adjacentList[0].push_back(3);
    adjacentList[0].push_back(4);
    adjacentList[0].push_back(5);
    adjacentList[1].push_back(0);
    adjacentList[1].push_back(2);
    adjacentList[1].push_back(3);
    adjacentList[2].push_back(0);
    adjacentList[2].push_back(1);
    adjacentList[3].push_back(0);
    adjacentList[3].push_back(1);
    adjacentList[4].push_back(0);
    adjacentList[4].push_back(5);
    adjacentList[5].push_back(0);
    adjacentList[5].push_back(4);
    bool isOld[6];
    for(int i=0; i<6; i++)
    {
        isOld[i] = false;
    }
    searchDepthFirst(0, adjacentList, isOld, 6, dfTree);
    for(vector<Edge>::iterator it = dfTree.begin(); it != dfTree.end(); it++)

```



```

    {
        cout << "(" << it -> start << ", " << it -> end << "): " << it -> weight << endl;
    }
}
int main()
{
    testSearchDepthFirst();
    system("pause");
    return 0;
}

```

A. 23 包含计算 LOW 值的深度优先搜索

```

#include <iostream>
#include <vector>
using namespace std;
struct Edge
{
    int start;
    int end;
    float weight;
    Edge(int _ start, int _ end, float _ weight)
    {
        start = _ start;
        end = _ end;
        weight = _ weight;
    }
    bool operator < (const Edge & e) const
    {
        return this -> weight < e.weight;
    }
    bool operator > (const Edge & e) const
    {
        return this -> weight > e.weight;
    }
    bool operator > = (const Edge & e) const
    {
        return this -> weight > = e.weight;
    }
    bool operator < = (const Edge & e) const
    {
        return this -> weight < = e.weight;
    }
    bool operator == (const Edge & e) const
    {
        return (this -> weight == e.weight &&
            this -> start == e.start &&
            this -> end == e.end);
    }
};

// Fig. 5.11 Depth-first search with LOW computation.
// 深度优先遍历求图的关节点的递归部分, 只被函数 searchDepthFirstLow() 调用
// 不被单独调用
void searchDepthFirstLowPrivate(int * dfnumber, int vertex, vector < int > * adjacentList,
int * low, int * father,
    bool * isOld, int nVertex, vector < int > & articulationPoints, vector <
Edge > & tree)
{
    // 检查输入参数是否有误

```

```

if(dfnumber == NULL ||
    adjacentList == NULL ||
    low == NULL ||
    isOld == NULL ||
    nVertex < 0)
{
    return;
}
static int count = 0;
int sonnum = 0;
isOld[vertex] = true;
dfnumber[vertex] = count;
count ++;
low[vertex] = dfnumber[vertex];
for(vector<int>::iterator it = adjacentList[vertex].begin();
    it != adjacentList[vertex].end();
    it ++)
{
    if(isOld[* it] == false)
    {
        tree.push_back(Edge(vertex, * it, 1));
        sonnum ++;
        father[* it] = vertex;
        searchDepthFirstLowPrivate(dfnumber, * it, adjacentList, low, father, isOld, nVertex, articulationPoints, tree);
        if(father[vertex] != -1 && low[* it] >= dfnumber[vertex]) // 不是根顶点
        {
            articulationPoints.push_back(vertex);
        }
        low[vertex] = min(low[vertex], low[* it]);
    }
    else if(* it != father[vertex])
    {
        low[vertex] = min(low[vertex], dfnumber[* it]);
    }
}
if(father[vertex] == -1 && sonnum > 1) // 根顶点, 且 son 的个数不只 1 个
{
    articulationPoints.push_back(vertex);
}
}
//通过图的深度优先遍历来找图的关节点
//用关节点来分隔各个连通分支
//输入:
//    vertex: 图中开始搜索的顶点
//    adjacentList: 数组, 存放所有顶点的邻接点列表
//    adjacentList[i]表示顶点 i 的所有邻接点的列表
//    nVertex: 图中顶点的数目, 顶点编号从 0 到 (nVertex - 1);
//    tree: 深度优先生成树的边集
//输出:
//    dfnumber: 图中的顶点在深度搜索中被遍历到的顺序序号
//    articulationPoints: 图中的关节点的列表
//测试函数:
//    void testSearchDepthFirstLow();
//测试数据:
//    遍历的起点是 v0
//    原图为如下:
//
//                (v3)
//                | \
//                |  \
//                |   \

```

```

//          |   \
//      (v4) - - (v0) - - (v1)
//          |   /   \   |
//          |   /     \   |
//          |   /       \   |
//      (v5)           (v2)
//
//测试结果:
//  生成的深度优先遍历树: (树的根是 v0)
//          (v3)
//          \
//          \
//          \
//      (v4) - - (v0) - - (v1)
//          |           |
//          |           |
//          |           |
//      (v5)           (v2)
//  得到的关节点是 v0
//第二组测试数据:
//  遍历的起点是 v0
//          v0
//          /   \
//          /     \
//          /       \
//      v1 - - - - v2
//          /   \
//          /     \
//      v3 - - - - v4
//          /
//          /
//      v5 - - - - v8
//          |   \   |
//          |   \   |
//          |   \   |
//      v6 - - - - v7
//第二组测试结果:
//  生成的深度遍历树为:
//          v0
//          /
//          /
//          /
//      v1 - - - - v2
//          /
//          /
//      v3 - - - - v4
//          /
//          /
//      v5      v8
//          |   |
//          |   |
//          |   |
//      v6 - - - - v7
//  得到的关节点为: v1, v3, v5
void searchDepthFirstLow(int * dfnumber, int vertex, vector<int> * adjacentList, int nVertex,
vector<int> & articulationPoints, vector<Edge> & tree)
{
    int * low = new int[nVertex];
    int * father = new int[nVertex];
    bool * isOld = new bool[nVertex];
    for(int i = 0; i < nVertex; i++)

```

```

    {
        low[i] = -1;
        father[i] = -1;
        dfnumber[i] = -1;
        isOld[i] = false;
    }
    tree.clear();
    searchDepthFirstLowPrivate(dfnumber, vertex, adjacentList, low, father, isOld, nVertex,
    articulationPoints, tree);
    if(low != NULL)
    {
        delete [] low;
    }
    if(father != NULL)
    {
        delete [] father;
    }
    if(isOld != NULL)
    {
        delete [] isOld;
    }
}
void testSearchDepthFirstLow()
{
    //第一组测试数据初始化
    vector<Edge> dfTree;
    //vector<int> adjacentList[6];
    //adjacentList[0].push_back(0);
    //adjacentList[0].push_back(1);
    //adjacentList[0].push_back(2);
    //adjacentList[0].push_back(3);
    //adjacentList[0].push_back(4);
    //adjacentList[0].push_back(5);
    //adjacentList[1].push_back(0);
    //adjacentList[1].push_back(2);
    //adjacentList[1].push_back(3);
    //adjacentList[2].push_back(0);
    //adjacentList[2].push_back(1);
    //adjacentList[3].push_back(0);
    //adjacentList[3].push_back(1);
    //adjacentList[4].push_back(0);
    //adjacentList[4].push_back(5);
    //adjacentList[5].push_back(0);
    //adjacentList[5].push_back(4);
    //int dfnumber[6];
    // 测试数据第二组
    vector<int> adjacentList[9];
    adjacentList[0].resize(2);
    adjacentList[0].at(0) = 1;
    adjacentList[0].at(1) = 2;
    adjacentList[1].resize(4);
    adjacentList[1].at(0) = 0;
    adjacentList[1].at(1) = 2;
    adjacentList[1].at(2) = 3;
    adjacentList[1].at(3) = 4;
    adjacentList[2].resize(2);
    adjacentList[2].at(0) = 0;
    adjacentList[2].at(1) = 1;
    adjacentList[3].resize(3);
    adjacentList[3].at(0) = 1;

```

```

adjacentList[3].at(1) = 4;
adjacentList[3].at(2) = 5;
adjacentList[4].resize(2);
adjacentList[4].at(0) = 1;
adjacentList[4].at(1) = 3;
adjacentList[5].resize(4);
adjacentList[5].at(0) = 3;
adjacentList[5].at(1) = 6;
adjacentList[5].at(2) = 7;
adjacentList[5].at(3) = 8;
adjacentList[6].resize(2);
adjacentList[6].at(0) = 5;
adjacentList[6].at(1) = 7;
adjacentList[7].resize(3);
adjacentList[7].at(0) = 5;
adjacentList[7].at(1) = 6;
adjacentList[7].at(2) = 8;
adjacentList[8].resize(2);
adjacentList[8].at(0) = 5;
adjacentList[8].at(1) = 7;
int dfnumber[9];
vector<int> articulationPoints;
articulationPoints.clear();
// 开始测试
//searchDepthFirstLow(dfnumber, 0, adjacentList, 6, articulationPoints, dfTree); //第一组
测试 searchDepthFirstLow(dfnumber, 0, adjacentList, 9, articulationPoints, dfTree); //第二组测
试
// 输出测试结果
for(vector<Edge>:: iterator it = dfTree.begin(); it != dfTree.end(); it++)
{
    cout << "(" << it->start << ", " << it->end << "): " << it->weight << endl;
}
cout << "关节点有: " << endl;
for(vector<int>:: iterator it = articulationPoints.begin(); it != articulation-
Points.end(); it++)
{
    cout << *it << '\t';
}
cout << endl;
}
int main()
{
    testSearchDepthFirstLow();
    system("pause");
    return 0;
}

```

A.24 计算有向图强连通分支的算法

```

#include<iostream>
#include<vector>
using namespace std;
bool isInStack(int elem, const vector<int>& stack)
{
    for(vector<int>:: const_ iterator it = stack.begin(); it != stack.end(); it++)
    {
        if(elem == *it) {
            return true;
        }
    }
}

```

```

    }
    return false;
}
// Fig 5.15 Procedure to compute LOWLINK.
// 深度优先遍历计算有向图强连通分支
// 输入:
//   vertex: 深度优先遍历起始点
//   isOld: 记录各点被访问的情况。isOld[i]为 true 表示点 i 已经被访问过了,
//          isOld[i]为 false 表示 i 尚未被访问过
//   adjacentList: 记录各点的邻接点链表
//   adjacentList[i]表示从点 i 链接到的点所组成的链表
//   nVertex: 图中所有点的个数
// 输出:
//   dfnumber: 深度优先遍历中点被遍历顺序序列
//   dfnumber[i]表示点 i 在深度优先遍历中被遍历的次序
//   lowLink: lowlink[i] 表示点 i 在以点 lowLink[i]为代表的连通分支内
// 测试数据:
//   原输入有向图如下:
//   点集: (0, 1, 2, 3, 4, 5, 6, 7);
//   边集: <0, 1> <0, 3> <0, 4>
//         <1, 2> <1, 3>
//         <2, 0>
//         <3, 2>
//         <4, 3>
//         <5, 6> <5, 7>
//         <6, 4>
//         <7, 3> <7, 5> <7, 6>
// 测试结果:
//   生成的强连通分支(0, 1, 2) (3)(4)(5, 7) (6)
void searchDepthFirstLowLink( int vertex, int * dfnumber, int * lowLink, bool * isOld, vector
<int> * adjacentList,
                             int nVertex, vector<int> & vertexStack)
{
    static int count;
    isOld[vertex] = true;
    dfnumber[vertex] = count;
    count ++;
    lowLink[vertex] = dfnumber[vertex];
    vertexStack.push_back(vertex);
    // 对于这个点的所有邻接点:
    for(vector<int>:: iterator it = adjacentList[vertex].begin();
        it != adjacentList[vertex].end();
        it ++)
    {
        if(isOld[* it] == false) // * it 这个点是新的, 尚未被遍历到
        {
            searchDepthFirstLowLink(* it, dfnumber, lowLink, isOld, adjacentList, nVertex,
vertexStack);
            lowLink[vertex] = min(lowLink[vertex], lowLink[* it]);
        }
        else if(dfnumber[* it] < dfnumber[vertex] &&
            isInStack(* it, vertexStack))
        {
            // 这个点已经被遍历过了
            lowLink[vertex] = min(dfnumber[* it], lowLink[vertex]);
        }
    }
} // end for-loop
if(lowLink[vertex] == dfnumber[vertex])
{

```

```
int x = -1;
do
{
    x = vertexStack.back();
    vertexStack.pop_back();
    cout << x << endl;
}while(x != vertex);
cout << "以上各点在同一下强连通分支!" << endl;
}
}

void testSearchDepthFirstLowLink()
{
    const int NVERTEX = 8;
    int dfnumber[NVERTEX];
    int lowlink[NVERTEX];
    bool isold[NVERTEX];
    vector<int> adjacentlist[NVERTEX];
    vector<int> vertexstack;
    // 初始化变量
    for(int i=0; i<NVERTEX; i++)
    {
        isold[i] = false;
        lowlink[i] = NVERTEX;
    }
    // 建有向图
    adjacentlist[0].push_back(1);
    adjacentlist[0].push_back(3);
    adjacentlist[0].push_back(4);
    adjacentlist[1].push_back(2);
    adjacentlist[1].push_back(3);
    adjacentlist[2].push_back(0);
    adjacentlist[4].push_back(3);
    adjacentlist[5].push_back(6);
    adjacentlist[5].push_back(7);
    adjacentlist[6].push_back(4);
    adjacentlist[7].push_back(3);
    adjacentlist[7].push_back(5);
    adjacentlist[7].push_back(6);
    // 在有向图中找生成树, 并找到强连通分支
    for(int i=0; i<NVERTEX; i++)
    {
        if(! isold[i]) // 找个点尚未被遍历到
        {
            searchDepthFirstLowLink(i, dfnumber, lowlink, isold, adjacentlist,
                                    NVERTEX, vertexstack);
        }
    }
    // 打印出 lowlink 结果
    for(int i=0; i<NVERTEX; i++)
    {
        cout << "dfnumber[" << i << "] = " << dfnumber[i] << " \t";
        cout << "lowlink[" << i << "] = " << lowlink[i] << endl;
    }
}

int main()
{
    testSearchDepthFirstLowLink();
    system("pause");
    return 0;
}
```

A.25 Dijkstra 算法

```

#include <iostream>
using namespace std;
// fig 5.24 Dijkstra's algorithm(shortest path).
// Dijkstra 计算从一点到多点的最短路径
// 输入:
//   startvertex: 作为起点得点
//   nVertex: 有向图的点数
//   weightTable: 有向图及其边权重的存放矩阵
//   weightTable[i * nVertex + j] 为从点 i 到点 j 的边的权重,
//   如果有向图中不存在该边, 则为一个最大数
// 输出:
//   dist: 从起点到各点的最短路径的长度
//   dist[i]表示从 startvertex 到点 i 的最短路径的长度
// 测试函数:
//   void testDijkstraShortestPath()
// 测试数据:
//   (2)
//   0 - - ->1
//   |   /   | (3)
// (10) |   / (7) |
//   \ / | /   \ /
//   4 - - ->2
//   / \ (6) /
//   |   / (4)
// (5) |   /
//   |   | /
//   3
// 测试结果:
//   0 ->1 最短路径长 2
//   0 ->2 最短路径长 5
//   0 ->3 最短路径长 9
//   0 ->4 最短路径长 9
int dijkstraShortestPath(int startvertex, int nVertex, int * weightTable, int * dist)
{
    if(startvertex >= nVertex)
    {
        cout << "Error: the start vertex is out of range. " << endl;
        return -1;
    }
    bool * isVisited = new bool[nVertex];

    // 临时数组的初始化
    for(int i=0; i<nVertex; i++)
    {
        // dist[i]中记录从 startvertex 到 i 的边的权重
        dist[i] = weightTable[nVertex * startvertex + i];
        isVisited[i] = false;
    }
    isVisited[startvertex] = true;
    bool allVisited = false;
    while(! allVisited)
    {
        int minVertex = -1;
        int minWeight = SHRT_MAX;
        for(int i=0; i<nVertex; i++)
        {
            if(isVisited[i] == false) // 点 i 不在 S 中, 在 V-S 中。
            {

```



```

        if (dist[i] < minWeight)
        {
            minVertex = i;
            minWeight = dist[i];
        }
    }

    if (minVertex == -1)
    {
        break;
    }
    isVisited[minVertex] = true;
    allVisited = true;
    for (int i = 0; i < nVertex; i++)
    {
        if (isVisited[i] == false) // 点 i 在 V-S 中
        {
            dist[i] = min(dist[i], dist[minVertex] + weightTable[nVertex * minVertex +
i]);

            allVisited = false;
        }
    }
}
// 释放临时数组空间
if (isVisited != NULL)
{
    delete [] isVisited;
    isVisited = NULL;
}
return 0;
}

void testDijkstraShortestPath()
{
    const int NVERTEX = 5;
    int weightTable[NVERTEX * NVERTEX];
    // 初始化带权有向图
    weightTable[0 * NVERTEX + 0] = 0;
    weightTable[0 * NVERTEX + 1] = 2;
    weightTable[0 * NVERTEX + 2] = SHRT_MAX;
    weightTable[0 * NVERTEX + 3] = SHRT_MAX;
    weightTable[0 * NVERTEX + 4] = 10;
    weightTable[1 * NVERTEX + 0] = SHRT_MAX;
    weightTable[1 * NVERTEX + 1] = 0;
    weightTable[1 * NVERTEX + 2] = 3;
    weightTable[1 * NVERTEX + 3] = SHRT_MAX;
    weightTable[1 * NVERTEX + 4] = 7;
    weightTable[2 * NVERTEX + 0] = SHRT_MAX;
    weightTable[2 * NVERTEX + 1] = SHRT_MAX;
    weightTable[2 * NVERTEX + 2] = 0;
    weightTable[2 * NVERTEX + 3] = 4;
    weightTable[2 * NVERTEX + 4] = SHRT_MAX;
    weightTable[3 * NVERTEX + 0] = SHRT_MAX;
    weightTable[3 * NVERTEX + 1] = SHRT_MAX;
    weightTable[3 * NVERTEX + 2] = SHRT_MAX;
    weightTable[3 * NVERTEX + 3] = 0;
    weightTable[3 * NVERTEX + 4] = 5;
    weightTable[4 * NVERTEX + 0] = SHRT_MAX;
    weightTable[4 * NVERTEX + 1] = SHRT_MAX;
    weightTable[4 * NVERTEX + 2] = 6;

```

```

weightTable[4 * NVERTEX + 3] = SHRT_MAX;
weightTable[4 * NVERTEX + 4] = 0;
int dist[NVERTEX];
dijkstraShortestPath(0, NVERTEX, (int *)weightTable, dist);
// 显示出从 0 到各个点的最短距离
for(int i=0; i<NVERTEX; i++)
{
    cout << "from 0 to " << i << " : " << dist[i] << endl;
}
}
int main()
{
    testDijkstraShortestPath();
    system("pause");
    return 0;
}

```

A.26 矩阵的 LUP 分解算法

```

#include <iostream>
#include <string>
#include <assert.h>
using namespace std;
/*****
    矩阵类:
    封装一部分矩阵基本操作
    *****/
class Matrix
{
public:
    Matrix(int row, int col)
    {
        int i=0;
        this->m_ppData=NULL;
        this->m_iRow=0;
        this->m_iCol=0;
        if(row<=0 || col<=0)
        {
            return;
        }
        this->m_iRow=row;
        this->m_iCol=col;
        this->m_ppData=new double* [row];
        for(i=0; i<row; i++)
        {
            this->m_ppData[i]=new double[col];
            memset(this->m_ppData[i], 0, sizeof(double) * this->m_iCol);
        }
    }
    Matrix(const Matrix& other)
    {
        this->operator=(other);
    }
    ~Matrix()
    {
        ClearMemory();
    }
    double* operator [](int rowIndex)
    {

```

```

double* pData = NULL;
if(rowIndex >= 0 && rowIndex < this->m_iRow)
{
    if(this->m_ppData != NULL)
    {
        pData = this->m_ppData[rowIndex];
    }
}
return pData;
}

Matrix& operator = (const Matrix& other)
{
    if(&other == this)
    {
        return *this;
    }
    int i=0;
    this->m_iCol = other.m_iCol;
    this->m_iRow = other.m_iRow;
    this->m_ppData = new double* [this->m_iRow];
    for(i=0; i < this->m_iRow; i++)
    {
        this->m_ppData[i] = new double[this->m_iCol];
        memcpy(this->m_ppData[i], other.m_ppData[i], sizeof(double) * this->m_
iCol);
    }
    return *this;
}

Matrix operator + (const Matrix& other)
{
    int i=0;
    int j=0;
    if(this->m_iCol != other.m_iCol || this->m_iRow != other.m_iRow)
    {
        return *this;
    }
    Matrix result(this->m_iRow, this->m_iCol);
    for(i=0; i < this->m_iRow; i++)
    {
        for(j=0; j < this->m_iCol; j++)
        {
            result[i][j] = (*this)[i][j] + const_cast<Matrix&>(other)[i][j];
        }
    }
    return result;
}

Matrix operator - (const Matrix& other)
{
    int i=0;
    int j=0;
    Matrix result(this->m_iRow, this->m_iCol);
    if(this->m_iCol != other.m_iCol || this->m_iRow != other.m_iRow)
    {
        return(*this);
    }
    for(i=0; i < this->m_iRow; i++)
    {
        for(j=0; j < this->m_iCol; j++)
        {

```

```

        result[i][j] = (* this)[i][j] - const_cast<Matrix &>(other)[i][j];
    }
}
return result;
}
friend ostream& operator << (ostream& o, const Matrix& a)
{
    if(a.m_ppData != NULL)
    {
        for(int i=0; i<a.m_iRow; i++)
        {
            for(int j=0; j<a.m_iCol; j++)
            {
                o<<a.m_ppData[i][j];
                o<<" ";
            }
            o<<endl;
        }
    }
    return o;
}
Matrix operator * (const Matrix& other)
{
    int i=0;
    int j=0;
    int k=0;
    Matrix result(this->m_iRow, other.m_iCol);
    if(this->m_iCol != other.m_iRow)
    {
        return result;
    }
    for(i=0; i<result.m_iRow; i++)
    {
        for(j=0; j<result.m_iCol; j++)
        {
            for(k=0; k<this->m_iCol; k++)
            {
                result[i][j] = result[i][j] +
                    (* this)[i][k] * const_cast<Matrix &>(other)[k][j];
            }
        }
    }
    return result;
}
void InterChangeColumn(int col1, int col2)
{
    if(col1 < 0 && col1 >= this->m_iCol)
    {
        return;
    }
    if(col2 < 0 && col2 >= this->m_iCol)
    {
        return;
    }
    for(int i=0; i<this->m_iRow; i++)
    {
        double temp = (* this)[i][col1];
        (* this)[i][col1] = (* this)[i][col2];
        (* this)[i][col2] = temp;
    }
}

```

```

    }
    void InitializeAsPermutationMatrix()
    {
        if(this->m_iCol != this->m_iRow)
        {
            return;
        }
        for(int i=0; i<this->m_iRow; i++)
        {
            (* this)[i][i]=1;
        }
    }
    int GetRowCount() { return m_iRow;}
    int GetColumnCount() { return m_iCol;}
private:
    void ClearMemory()
    {
        int i=0;
        if(this->m_ppData != NULL)
        {
            for(i=0; i<this->m_iRow; i++)
            {
                delete [] (this->m_ppData[i]);
                this->m_ppData[i]=NULL;
            }
            delete [] this->m_ppData;
            this->m_ppData=NULL;
        }
    }
private:
    double* * m_ppData;
    int m_iRow;
    int m_iCol;
};

Matrix GetUpperTriangularInverse(Matrix& a);
// LUP 矩阵分解算法
// 功能: 将一 n × n 矩阵分解为单位下三角矩阵, 上三角矩阵和组合矩阵的乘积
// 函数: void FACTOR(Matrix& A, int m, int p, Matrix* & L, Matrix* & U, Matrix* & P)
// 输入:
//       Matrix& A: 矩阵 A
//       int m: 矩阵 A 的行数
//       int p: 矩阵 A 的列数
// 输出:
//       Matrix* & L: 分解出的单位下三角矩阵
//       Matrix* & U: 分解出的上三角矩阵
//       Matrix* & P: 分解出的组合矩阵
// 测试函数: TestFACTOR()
// 输入:
//       A:      0   0   0   1
//              0   0   2   0
//              0   3   0   0
//              4   0   0   0
//
//       m:      4
//       p:      4
// 输出:
//       L:      1   0   0   0
//              0   1   0   0
//              0   0   1   0
//              0   0   0   1

```

```

//
//      U:      1   0   0   0
//              0   2   0   0
//              0   0   3   0
//              0   0   0   4
//
//      P:      0   0   0   1
//              0   0   1   0
//              0   1   0   0
//              1   0   0   0
void FACTOR(Matrix& A, int m, int p, Matrix* & L, Matrix* & U, Matrix* & P)
{
    if(m == 1)
    {
        L = new Matrix(1, 1);
        (* L)[0][0] = 1;

        // 将 P 初始化为单位矩阵
        P = new Matrix(p, p);
        P->InitializeAsPermutationMatrix();
        // 找到 A 的有非 0 元素的第 c 行
        int iARow = A.GetRowCount();
        int iACol = A.GetColumnCount();
        int c = 0;
        bool found = false;
        for(c = 0; c < iACol; c++)
        {
            for(int r = 0; r < iARow; r++)
            {
                if(A[r][c] != 0)
                {
                    found = true;
                    break;
                }
            }
            if(found)
            {
                break;
            }
        }
        if(found)
        {
            P->InterChangeColumn(0, c);
        }

        U = new Matrix(iARow, iACol);
        * U = A * (* P);
        return;
    }
    else
    {
        Matrix B(m / 2, p);
        Matrix C(m / 2, p);
        int i = 0;
        int j = 0;
        for(i = 0; i < m / 2; i++)
        {
            for(j = 0; j < p; j++)

```

```

    {
        B[i][j] = A[i][j];
    }
}
for( ; i < m; i++)
{
    for(j=0; j < p; j++)
    {
        C[i-m/2][j] = A[i][j];
    }
}

Matrix* L1 = NULL;
Matrix* U1 = NULL;
Matrix* P1 = NULL;
FACTOR(B, m / 2, p, L1, U1, P1);

Matrix D = C * (* P1); // P = P 的逆

Matrix E(m/2, m/2);
Matrix F(m/2, m/2);
for(i=0; i < m/2; i++)
{
    for(j=0; j < m/2; j++)
    {
        E[i][j] = (* U1)[i][j];
    }
}
for(i=0; i < m/2; i++)
{
    for(j=0; j < m/2; j++)
    {
        F[i][j] = D[i][j];
    }
}

Matrix EInverse(m / 2, m / 2);
// 计算 E 的逆矩阵
// 因为 P1 是上三角矩阵, 可用特殊方法求它的逆
EInverse = GetUpperTriangularInverse(E);
Matrix G = D - F * EInverse * (* U1);

Matrix Gother(m / 2, p - m / 2);
for(i=0; i < m/2; i++)
{
    for(j=0; j < p - m / 2; j++)
    {
        Gother[i][j] = G[i][j + m / 2];
    }
}

Matrix * L2 = NULL;
Matrix* U2 = NULL;
Matrix* P2 = NULL;
FACTOR(Gother, m / 2, p - m / 2, L2, U2, P2);

Matrix P3(p, p);
for(i=0; i < m/2; i++)

```

```

{
    P3[i][i] = 1;
}
for(i = m / 2; i < p; i++)
{
    for(j = m / 2; j < p; j++)
    {
        P3[i][j] = (* P2)[i - m / 2][j - m / 2];
    }
}
Matrix H = (* U1) * (P3);
L = new Matrix(m, m);
for(i = 0; i < m / 2; i++)
{
    for(j = 0; j < m / 2; j++)
    {
        (* L)[i][j] = (* L1)[i][j];
    }
}

Matrix temp = F * EInverse;
for(i = m / 2; i < m; i++)
{
    for(j = 0; j < m / 2; j++)
    {
        (* L)[i][j] = temp[i - m / 2][j];
    }
}

for(i = m / 2; i < m; i++)
{
    for(j = m / 2; j < m; j++)
    {
        (* L)[i][j] = (* L2)[i - m / 2][j - m / 2];
    }
}
U = new Matrix(m, p);
for(i = 0; i < m / 2; i++)
{
    for(j = 0; j < p; j++)
    {
        (* U)[i][j] = H[i][j];
    }
}

for(i = m / 2; i < m; i++)
{
    for(j = m / 2; j < p; j++)
    {
        (* U)[i][j] = (* U2)[i - m / 2][j - m / 2];
    }
}

P = new Matrix(p, p);
(* P) = P3 * (* P1);
delete L1;
delete U1;
delete P1;
delete L2;
delete U2;

```



```

        delete P2;
    }
}
/* *****
    求上三角矩阵的逆矩阵
    ***** */
Matrix GetUpperTriangularInverse(Matrix& a)
{
    int r = a.GetRowCount();
    int c = a.GetColumnCount();
    assert(r == c);
    Matrix res(r, c);
    for(int i=0; i<r; i++)
    {
        res[i][i] = 1.0 / a[i][i];
    }
    for(int i=0; i<r; i++)
    {
        for(int j=0; j<c; j++)
        {
            if(i < j)
            {
                double t=0;
                for(int k=i; k<j-1; k++)
                {
                    t += a[j][k] * res[k][i];
                }
                res[i][j] = -1 * (1.0 / a[i][i]) * t;
            }
        }
    }
    return res;
}
//
//    打印输出矩阵
//
void PrintMatrix(string name, Matrix& m)
{
    cout << " - - - - - " << name << " - - - - - " << endl;
    cout << endl;
    cout << m;
}
void TestFACTOR()
{
    Matrix a(4, 4);
    a[0][3] = 1;
    a[1][2] = 2;
    a[2][1] = 3;
    a[3][0] = 4;
    Matrix * L=NULL;
    Matrix* U=NULL;
    Matrix* P=NULL;
    FACTOR(a, 4, 4, L, U, P);
    PrintMatrix("L", * L);
    PrintMatrix("U", * U);
    PrintMatrix("P", * P);
}
int main(int argc, char* argv[])
{
    TestFACTOR();
}

```

```

    system("pause");
    return 0;
}

```

A.27 Four Russians 布尔矩阵相乘算法

```

#include <iostream>
#include <string>
#include <math.h>
#include <assert.h>
using namespace std;
// FourRussians 布尔矩阵乘法
// 功能: 求两个 N * N 的布尔矩阵的乘积
// 函数: void FourRussiansBooleanMatrixMultiplication(Matrix& A, Matrix& B, Matrix& C)
// 输入:
//       Matrix& A: 矩阵 A
//       Matrix& B: 矩阵 B
// 输出:
//       Matrix& C: 矩阵 C
// 测试函数: TestFourRussian()
// 测试数据:
//
//       ----- A -----
//
//       0      0      1      0      0      0      0      0
//       1      0      1      0      0      0      0      0
//       1      1      1      0      0      0      0      0
//       1      0      0      0      0      0      0      0
//       0      0      0      0      0      0      0      0
//       1      1      0      0      0      0      0      0
//       0      0      0      0      0      0      0      0
//       0      1      1      0      0      0      0      0
//
//       ----- B -----
//
//       0      1      0      1      1      0      0      1
//       0      0      0      1      0      1      0      0
//       1      1      0      1      0      0      0      0
//       0      0      0      0      0      0      0      0
//       0      0      0      0      0      0      0      0
//       0      0      0      0      0      0      0      0
//       0      0      0      0      0      0      0      0
//       0      0      0      0      0      0      0      0
//
// 测试函数:
//
//       ----- C -----
//
//       1      1      0      1      0      0      0      0
//       1      1      0      1      1      0      0      1
//       1      1      0      1      1      1      0      1
//       0      1      0      1      1      0      0      1
//       0      0      0      0      0      0      0      0
//       0      1      0      1      1      1      0      1
//       0      0      0      0      0      0      0      0
//       1      1      0      1      0      1      0      0
//
//
// 计算以 2 为底的 log
//
double log2(double a)
{
    return log(a) / log(2.0);
}

```

```

/* *****
   矩阵类:
   封装一部分矩阵基本操作
   ***** */
class Matrix
{
public:
    Matrix(int row, int col)
    {
        int i=0;
        this->m_ppData=NULL;
        this->m_iRow=0;
        this->m_iCol=0;
        if(row<=0 || col<=0)
        {
            return;
        }
        this->m_iRow=row;
        this->m_iCol=col;
        this->m_ppData=new double* [row];
        for(i=0; i<row; i++)
        {
            this->m_ppData[i]=new double[col];
            memset(this->m_ppData[i], 0, sizeof(double) * this->m_iCol);
        }
    }
    Matrix(const Matrix& other)
    {
        this->operator=(other);
    }
    ~Matrix()
    {
        ClearMemory();
    }
    double* operator [] (int rowIndex)
    {
        double* pData=NULL;
        if(rowIndex>=0 && rowIndex<this->m_iRow)
        {
            if(this->m_ppData!=NULL)
            {
                pData=this->m_ppData[rowIndex];
            }
        }
        return pData;
    }
    Matrix& operator=(const Matrix& other)
    {
        if(&other==this)
        {
            return *this;
        }
        int i=0;
        this->m_iCol=other.m_iCol;
        this->m_iRow=other.m_iRow;
        this->m_ppData=new double* [this->m_iRow];
        for(i=0; i<this->m_iRow; i++)
        {
            this->m_ppData[i]=new double[this->m_iCol];
            memcpy(this->m_ppData[i], other.m_ppData[i], sizeof(double) * this->m_iCol);
        }
    }
};

```

```

iCol);
    }
    return * this;
}
Matrix operator + (const Matrix& other)
{
    int i=0;
    int j=0;
    if(this->m_iCol != other.m_iCol || this->m_iRow != other.m_iRow)
    {
        return * this;
    }
    Matrix result(this->m_iRow, this->m_iCol);
    for(i=0; i<this->m_iRow; i++)
    {
        for(j=0; j<this->m_iCol; j++)
        {
            result[i][j] = (* this)[i][j] + const_cast<Matrix&>(other)[i][j];
        }
    }
    return result;
}
Matrix operator - (const Matrix& other)
{
    int i=0;
    int j=0;
    Matrix result(this->m_iRow, this->m_iCol);
    if(this->m_iCol != other.m_iCol || this->m_iRow != other.m_iRow)
    {
        return(* this);
    }
    for(i=0; i<this->m_iRow; i++)
    {
        for(j=0; j<this->m_iCol; j++)
        {
            result[i][j] = (* this)[i][j] - const_cast<Matrix&>(other)[i][j];
        }
    }
    return result;
}
friend ostream& operator << (ostream& o, const Matrix& a)
{
    if(a.m_ppData != NULL)
    {
        for(int i=0; i<a.m_iRow; i++)
        {
            for(int j=0; j<a.m_iCol; j++)
            {
                o<<a.m_ppData[i][j];
                o<<" ";
            }
            o<<endl;
        }
    }
    return o;
}
Matrix operator * (const Matrix& other)
{
    int i=0;

```

```

int j=0;
int k=0;
Matrix result(this->m_iRow, other.m_iCol);
if(this->m_iCol != other.m_iRow)
{
    return result;
}
for(i=0; i<result.m_iRow; i++)
{
    for(j=0; j<result.m_iCol; j++)
    {
        for(k=0; k<this->m_iCol; k++)
        {
            result[i][j] = result[i][j] +
                (* this)[i][k] * const_cast<Matrix &>(other)[k][j];
        }
    }
}
return result;
}

void InterChangeColumn(int col1, int col2)
{
    if(col1 < 0 && col1 >= this->m_iCol)
    {
        return;
    }
    if(col2 < 0 && col2 >= this->m_iCol)
    {
        return;
    }
    for(int i=0; i<this->m_iRow; i++)
    {
        double temp = (* this)[i][col1];
        (* this)[i][col1] = (* this)[i][col2];
        (* this)[i][col2] = temp;
    }
}

void InitializeAsPermutationMatrix()
{
    if(this->m_iCol != this->m_iRow)
    {
        return;
    }
    for(int i=0; i<this->m_iRow; i++)
    {
        (* this)[i][i] = 1;
    }
}

int GetRowCount() { return m_iRow;}
int GetColumnCount() { return m_iCol;}

private:
void ClearMemory()
{
    int i=0;
    if(this->m_ppData != NULL)
    {
        for(i=0; i<this->m_iRow; i++)
        {
            delete [] (this->m_ppData[i]);
            this->m_ppData[i] = NULL;
        }
    }
}

```

```

    }
    delete []this->m_ppData;
    this->m_ppData = NULL;
}

private:
    double* * m_ppData;
    int m_iRow;
    int m_iCol;
};
//
//      打印输出矩阵
//
void PrintMatrix(string name, Matrix& m)
{
    cout << "-----" << name << "-----" << endl;
    cout << endl;
    cout << m;
}
//
// 得到由 0 和 1 组成的逆序的整数
// 例如 NUM([0, 1, 1]) = 6
//
long NUM(Matrix& a)
{
    assert(a.GetRowCount() == 1);
    int cols = a.GetColumnCount();
    long sum = 0;
    for(int i = 0; i < cols; i++)
    {
        sum += (long)(a[0][i]) * (long)pow(2.0, i);
    }
    return sum;
}

void FourRussiansBooleanMatrixMultipulation(Matrix& A, Matrix& B, Matrix& C)
{
    assert(A.GetRowCount() == A.GetColumnCount() && B.GetRowCount() == B.GetColumnCount() &&
    A.GetRowCount() == B.GetRowCount());
    int n = A.GetRowCount();
    int m = (int)log2((double)n);
    int numberOfPartition = (int)(n / (double)m + 0.5);
    // 将 A 和 B 分到不同矩阵
    Matrix* * APartitions = new Matrix* [numberOfPartition];
    for(int i = 0; i < numberOfPartition - 1; i++)
    {
        APartitions[i] = new Matrix(n, m);
        for(int j = 0; j < n; j++)
        {
            for(int k = 0; k < m; k++)
            {
                (*APartitions[i])[j][k] = A[j][i * m + k];
            }
        }
    }
    APartitions[numberOfPartition - 1] = new Matrix(n, m);
    for(int j = 0; j < n; j++)
    {
        for(int k = 0; k < m; k++)
        {
            if(k < n - m * (numberOfPartition - 1))

```

```

        {
            (* APartitions[numberOfPartition-1])[j][k] = A[j][(numberOfPartition-1) *
m+k];
        }
        else
        {
            (* APartitions[numberOfPartition-1])[j][k] = 0;
        }
    }
}

Matrix* * BPartitions = new Matrix* [numberOfPartition];
for(int i=0; i<numberOfPartition-1; i++)
{
    BPartitions[i] = new Matrix(m, n);
    for(int j=0; j<m; j++)
    {
        for(int k=0; k<n; k++)
        {
            (* BPartitions[i])[j][k] = B[i * m+j][k];
        }
    }
}
BPartitions[numberOfPartition-1] = new Matrix(m, n);
for(int j=0; j<m; j++)
{
    for(int k=0; k<n; k++)
    {
        if((numberOfPartition-1) * m+j<n)
        {
            (* BPartitions[numberOfPartition-1])[j][k] = B[(numberOfPartition-1) * m+
j][k];
        }
        else
        {
            (* BPartitions[numberOfPartition-1])[j][k] = 0;
        }
    }
}

Matrix* * CPartitions = new Matrix* [numberOfPartition];

for(int i=0; i<numberOfPartition; i++)
{
    Matrix* * ROWSUM = new Matrix* [(long)pow(2.0, m)];

    ROWSUM[0] = new Matrix(1, n);

    for(int j=1; j<(long)pow(2.0, m); j++)
    {
        int k=0;
        while(true)
        {
            long value1 = (long)pow(2.0, k);
            long value2 = (long)pow(2.0, k+1);
            if(value1 <= j && j < value2)
            {
                break;
            }
        }
    }
}

```

```

        else
        {
            k++;
        }
    }

    ROWSUM[j] = new Matrix(1, n);
    Matrix bRow(1, n);
    for(int index=0; index<n; index++)
    {
        bRow[0][index] = (* BPartitions[i])[k][index];
    }
    (* ROWSUM[j]) = (* ROWSUM[j - (long)pow(2.0, k)]) + bRow;
    for(int index=0; index<n; index++)
    {
        (* ROWSUM[j])[0][index] = (* ROWSUM[j])[0][index] > 0 ? 1 : 0;
    }
}

CPartitions[i] = new Matrix(n, n);
for(int j=0; j<n; j++)
{
    Matrix aRow(1, m);
    for(int index=0; index<m; index++)
    {
        aRow[0][index] = (* APartitions[i])[j][index];
    }
    long numA = NUM(aRow);

    for(int index=0; index<n; index++)
    {
        (* CPartitions[i])[j][index] = (* ROWSUM[numA])[0][index];
    }
}

for(int index=0; index<(long)pow(2.0, m); index++)
{
    delete ROWSUM[index];
    ROWSUM[index] = NULL;
}
delete ROWSUM;
}

for(int index=0; index<numberOfPartition; index++)
{
    Matrix temp = C + (* CPartitions[index]);
    C = C + temp;
}

for(int i=0; i<n; i++)
{
    for(int j=0; j<n; j++)
    {
        C[i][j] = C[i][j] > 0 ? 1 : 0;
    }
}

for(int i=0; i<numberOfPartition; i++)
{

```



```

        delete APartitions[i];
        APartitions[i] = NULL;
    }
    delete APartitions;
    for(int i=0; i < numberOfPartition; i++)
    {
        delete BPartitions[i];
        BPartitions[i] = NULL;
    }
    delete BPartitions;

    for(int i=0; i < numberOfPartition; i++)
    {
        delete CPartitions[i];
        CPartitions[i] = NULL;
    }
    delete CPartitions;
}
void TestFourRussian()
{
    int n=8;
    Matrix A(n, n);
    Matrix B(n, n);
    Matrix C(n, n);
    A[0][2] = 1;
    A[1][0] = 1;
    A[1][2] = 1;
    A[2][0] = 1;
    A[2][1] = 1;
    A[2][2] = 1;
    A[3][0] = 1;
    A[5][0] = 1;
    A[5][1] = 1;
    A[7][1] = 1;
    A[7][2] = 1;
    B[0][1] = 1;
    B[0][3] = 1;
    B[0][4] = 1;
    B[0][7] = 1;
    B[1][3] = 1;
    B[1][5] = 1;
    B[2][0] = 1;
    B[2][1] = 1;
    B[2][3] = 1;
    FourRussiansBooleanMatrixMultipulation(A, B, C);
    PrintMatrix("A", A);
    PrintMatrix("B", B);
    PrintMatrix("C", C);
}
int main()
{
    TestFourRussian();
    system("pause");
    return 0;
}

```

A. 28 快速傅里叶变换算法

```

#include <iostream>
#include <vector>

```

```

#include <complex>
using namespace std;
/*****
 *   FFT
 *   功能: 快速傅里叶变换的实现
 *   函数: void FFT(vector<COMPLEX>a, vector<COMPLEX>&b)
 *   输入:
 *           vector<COMPLEX>a: 用系数来表示的多项式
 *   输出:
 *           vector<COMPLEX>&b: 复数表示的傅里叶变换的结果
 *   测试函数: TestFFT()
 *   输入:
 *           a: 多项式  $0 + x + 2x^2 + 3x^3 + 4x^4 + 5x^5 + 6x^6 + 7x^7$ , 表示为[0, 1, 2, 3, 4, 5, 6, 7]
 *   输出:
 *           b: [28, -4 - i* 9.65685, -4 - i* 4, -4 - i* 1.65685, -4, -4 + i* 1.65685, -4 + i* 4,
-4 + i* 9.65685]
 *
 *****/
typedef vector<unsigned char> Bits;
typedef complex<double> COMPLEX;
const double PI = 3.1415926535897931;
const double HalfPI = PI / 2;
//
// 计算以 2 为底的 log
//
double log2(double a)
{
    return log(a) / log(2.0);
}
//
// 将一个整数表示为二进制 0, 1 格式
// 如 3 表示为[1, 1], 6 表示为[1, 0, 0]
// i 代表要表示为二进制的整数, k 代表二进制数组的位数
//
Bits ToBinary(unsigned int i, int k)
{
    Bits bits;
    do
    {
        bits.insert(bits.begin(), (i % 2));
        i /= 2;
    } while(i != 0);
    while(bits.size() != k)
    {
        bits.insert(bits.begin(), 0);
    }
    return bits;
}
//
// 将一个二进制表示的数还原成整数形式
//
unsigned int ToInteger(Bits& bits)
{
    unsigned result = 0;
    for(unsigned int i = 0; i < bits.size(); i++)
    {
        result = result * 2 + bits[i];
    }
    return result;
}

```

```

//
// 假设整数 j 用二进制表示为 [d0, d1, d2, d3... dk-1]
// 那么本函数求二进制表示的 [dk-1, dk-2, ... d3, d2, d1, d0] 的值
//
unsigned int rev(unsigned int j, unsigned int k)
{
    Bits bits = ToBinary(j, k);

    for(unsigned i=0; i<k/2; i++)
    {
        unsigned char temp = bits[i];
        bits[i] = bits[k-1-i];
        bits[k-1-i] = temp;
    }
    return ToInteger(bits);
}
// 用来求 -1 的单位根
COMPLEX w(int n, int power)
{
    power = power % n;
    double u = 2 * PI * power / (double)n;
    // 判断几个特殊值
    double threshold = 0.000000001;
    if(abs(u-0) < threshold)
    {
        return COMPLEX(1, 0);
    }
    else if(abs(u - HalfPI) < threshold)
    {
        return COMPLEX(0, 1);
    }
    else if(abs(u - PI) < threshold)
    {
        return COMPLEX(-1, 0);
    }
    else if(abs(u - 3 * HalfPI) < threshold)
    {
        return COMPLEX(0, -1);
    }
    else if(abs(u - 2 * PI) < threshold)
    {
        return COMPLEX(1, 0);
    }
    else
    {
        return COMPLEX(cos(u), sin(u));
    }
}
}
void FFT(vector<COMPLEX> a, vector<COMPLEX> &b)
{
    // a 的长度必须是 2 个 n 次方
    int n = (int)a.size();
    int k = (int)log2(n);
    int lMax = n;
    int mMax = k+1;
    // 因为 l=n 并且 m=k+1
    // 所以要分配 1 * m 的空间

    vector<COMPLEX> * * * pppRPolynomial = new vector<COMPLEX> * * [lMax];
    for(int index=0; index<lMax; index++)

```

```

{
    pppRPolynomial[index] = new vector <COMPLEX> * [mMax];
    for(int mIndex=0; mIndex<mMax; mIndex++)
    {
        pppRPolynomial[index][mIndex] = new vector <COMPLEX> ();
    }
}

(* pppRPolynomial[0][k]) = a;
for(int m=k-1; m>=0; m--)
{
    for(int l=0; l<n; l += (int)(pow(2.0, m+1)))
    {
        a = (* pppRPolynomial[l][m+1]);

        unsigned int s = rev((unsigned int)(1/(int)pow(2.0, m)), k);

        pppRPolynomial[l][m] -> resize((unsigned int)pow(2.0, m));
        for(unsigned int j=0; j< (unsigned int)pow(2.0, m); j++)
        {
            (* pppRPolynomial[l][m])[j] = (a[j] + a[j + (unsigned int)pow(2.0, m)] * w(n,
s));
        }
        pppRPolynomial[l + (unsigned int)pow(2.0, m)][m] -> resize((unsigned int)pow(2.0,
m));
        for(unsigned int j=0; j< (unsigned int)pow(2.0, m); j++)
        {
            (* pppRPolynomial[l + (unsigned int)pow(2.0, m)][m])[j] = (a[j] + a[j + (un-
signed int)pow(2.0, m)] * w(n, s + n / 2));
        }
    }
}

// 输出结果
b.resize(n);
for(int l=0; l<n; l++)
{
    b[rev(l, k)] = (* pppRPolynomial[l][0])[0];
}

// 释放内存
for(int i=0; i<lMax; i++)
{
    for(int j=0; j<mMax; j++)
    {
        if(pppRPolynomial[i][j] != NULL)
        {
            delete pppRPolynomial[i][j];
        }
    }
    delete []pppRPolynomial[i];
}
delete []pppRPolynomial;
}

void TestFFT()
{
    int n=8;
    vector <COMPLEX> a;
    vector <COMPLEX> b;
    for(int i=0; i<n; i++)
    {
        a.push_back(COMPLEX(i, 0));
    }
}

```

```

    }
    FFT(a, b);
    cout << "结果是:" << endl;
    for(vector<COMPLEX>:: iterator it = b.begin();
        it != b.end();
        it++)
    {
        cout << it->real() << " + i * " << it->imag() << endl;
    }
}

int main()
{
    TestFFT();
    system("pause");
    return 0;
}

```

A.29 简化的傅里叶变换算法

```

#include <iostream>
#include <vector>
#include <complex>
using namespace std;
// FFT
// 功能: 快速傅里叶变换的实现
// 函数: void FFT(vector<COMPLEX> a, vector<COMPLEX> &b)
// 输入:
//       vector<COMPLEX> a: 用系数来表示的多项式
// 输出:
//       vector<COMPLEX> &b: 复数表示的傅里叶变换的结果
// 测试函数: TestFFT()
// 输入:
// a: 多项式  $0 + x + 2x^2 + 3x^3 + 4x^4 + 5x^5 + 6x^6 + 7x^7$ , 表示为 [0, 1, 2, 3, 4, 5, 6, 7]
// 输出:
// b: [28, -4 - i* 9.65685, -4 - i* 4, -4 - i* 1.65685, -4, -4 + i* 1.65685, -4 + i* 4,
-4 + i* 9.65685]
//
typedef vector<unsigned char> Bits;
typedef complex<double> COMPLEX;
const double PI = 3.1415926535897931;
const double HalfPI = PI / 2;
//
// 计算以 2 为底的 log
//
double log2(double a)
{
    return log(a) / log(2.0);
}
//
// 将一个整数表示为二进制 0, 1 格式
// 如 3 表示为 [1, 1], 6 表示为 [1, 0, 0]
// i 代表要表示为二进制的整数, k 代表二进制数组的位数
//
Bits ToBinary(unsigned int i, int k)
{
    Bits bits;
    do
    {
        bits.insert(bits.begin(), (i % 2));
        i /= 2;
    }
    while (i > 0);
    while (bits.size() < k)
        bits.insert(bits.begin(), 0);
}

```

```

    }while(i != 0);
    while(bits.size() != k)
    {
        bits.insert(bits.begin(), 0);
    }
    return bits;
}
//
//      将一个二进制表示的数还原成整数形式
//
unsigned int ToInteger(Bits& bits)
{
    unsigned result = 0;
    for(unsigned int i = 0; i < bits.size(); i++)
    {
        result = result * 2 + bits[i];
    }
    return result;
}
//
//      假设整数 j 用二进制表示为 [d0, d1, d2, d3... dk-1]
//      那么本函数求二进制表示的 [dk-1, dk-2, ... d3, d2, d1, d0] 的值
//
unsigned int rev(unsigned int j, unsigned int k)
{
    Bits bits = ToBinary(j, k);

    for(unsigned i = 0; i < k / 2; i++)
    {
        unsigned char temp = bits[i];
        bits[i] = bits[k - 1 - i];
        bits[k - 1 - i] = temp;
    }
    return ToInteger(bits);
}
// 用来求 -1 的单位根
COMPLEX w(int n, int power)
{
    power = power % n;
    double u = 2 * PI * power / (double)n;
    // 判断几个特殊值
    double threshold = 0.0000000001;
    if(abs(u - 0) < threshold)
    {
        return COMPLEX(1, 0);
    }
    else if(abs(u - HalfPI) < threshold)
    {
        return COMPLEX(0, 1);
    }
    else if(abs(u - PI) < threshold)
    {
        return COMPLEX(-1, 0);
    }
    else if(abs(u - 3 * HalfPI) < threshold)
    {
        return COMPLEX(0, -1);
    }
    else if(abs(u - 2 * PI) < threshold)
    {

```

```

        return COMPLEX(1, 0);
    }
    else
    {
        return COMPLEX(cos(u), sin(u));
    }
}

void SimplifiedFFT(vector<COMPLEX> a, vector<COMPLEX> &b)
{
    // a 的长度必须是 2 个 n 次方
    int n = (int)a.size();
    int k = (int)log2(n);
    int lMax = n;
    int mMax = k + 1;

    vector<COMPLEX> R;
    vector<COMPLEX> S;
    R = a;
    for(int l = 0; l < k; l++)
    {
        S = R;
        for(int i = 0; i < n; i++)
        {
            Bits bits = ToBinary(i, k);
            Bits SIndexBits = bits;
            SIndexBits[l] = 0;
            Bits wIndexBits;
            for(int index = 1; index >= 0; index--)
            {
                wIndexBits.push_back(bits[index]);
            }
            for(int index = l + 1; index < (int)bits.size(); index++)
            {
                wIndexBits.push_back(0);
            }
            Bits SIndexBits2 = bits;
            SIndexBits2[l] = 1;

            int SIndex = ToInteger(SIndexBits);
            int wIndex = ToInteger(wIndexBits);
            int SIndex2 = ToInteger(SIndexBits2);
            R[i] = S[SIndex] + w(n, wIndex) * S[SIndex2];
        }
    }
    b.resize(n);
    for(int i = 0; i < n; i++)
    {
        b[i] = R[rev(i, k)];
    }
}

void TestSimplifiedFFT()
{
    int n = 8;
    vector<COMPLEX> a;
    vector<COMPLEX> b;
    for(int i = 0; i < n; i++)
    {
        a.push_back(COMPLEX(i, 0));
    }
    SimplifiedFFT(a, b);
}

```

```

for(vector<COMPLEX>:: iterator it = b.begin();
    it != b.end();
    it++)
{
    cout << it->real() << " + i * " << it->imag() << endl;
}
}
int main()
{
    TestSimplifiedFFT();
    system("pause");
    return 0;
}

```

A. 30 求整数倒数的算法

```

#include <iostream>
#include <vector>
#include <assert.h>
using namespace std;
// RECIPROCAL
// 功能: 求整数倒数
// 函数: void RECIPROCAL(BitArray& P, BitArray& Value)
// 输入:
//      BitArray& P: BIT 方式表示的整数
// 输出:
//      BitArray& Value: BIT 方式表示的 P 的倒数
// 测试函数: TestRECIPROCAL()
// 输入:
//      P: [10011001]
// 输出:
//      Value: [011010110]
//
// 为了简化操作, 使用 unsigned char 来表示一个 bit
// bit 数组: 将一个整数表示为一系列的 bit 值 如 2 进制数 [10111] 表示为 1, 1, 1, 0, 1
//
class BitArray : public vector<unsigned char>
{
public:
    BitArray& operator = (const BitArray& other)
    {
        this->clear();
        this->resize(other.size());
        for(unsigned int i=0; i<other.size(); i++)
        {
            (*this)[i] = other[i];
        }
        return *this;
    }
    BitArray operator* (const BitArray& other)
    {
        BitArray result;
        for(unsigned i=0; i<other.size(); i++)
        {
            BitArray temp;
            for(unsigned j=0; j<i; j++)
            {
                temp.insert(temp.begin(), 0);
            }

```



```

        for(unsigned j=0; j<this->size(); j++)
        {
            temp.push_back(other[i] * (*this)[j]);
        }

        temp.RemoveZero();
        result=result+temp;

        result.Smooth();

    }
    result.Smooth();
    return result;
}

BitArray operator + (const BitArray& other)
{
    BitArray a;
    unsigned i=0;
    for( i=0; i<this->size() && i<other.size(); i++)
    {
        a.push_back((*this)[i]+other[i]);
    }
    if(i<this->size())
    {
        for(; i<this->size(); i++)
        {
            a.push_back((*this)[i]);
        }
    }
    else if(i<other.size())
    {
        for(; i<other.size(); i++)
        {
            a.push_back(other[i]);
        }
    }
    a.Smooth();
    return a;
}

bool operator <= (const BitArray& other)
{
    if(this->size()<other.size())
    {
        return true;
    }
    if(this->size()>other.size())
    {
        return false;
    }
    BitArray result;
    for(unsigned i=0; i<other.size(); i++)
    {
        if((*this)[i]>=other[i])
        {
            result.push_back((*this)[i]-other[i]);
        }
        else
        {
            if((*this)[i+1]==0)
            {

```

```

        return true;
    }
    else
    {
        (* this)[i+1] -= 1;
        result.push_back(1);
    }
}

}
result.RemoveZero();
if(result.size() == 0)
{
    return true;
}
else
{
    return false;
}
}

BitArray operator - (const BitArray& other)
{
    // 假设 this 的值比 other 要大
    assert(this->size() >= other.size());
    if(this->size() == other.size())
    {
        assert((* this)[this->size()-1] >= other[other.size()-1]);
    }
    BitArray result;
    BitArray tempThis = * this;
    unsigned int i = 0;
    for(i=0; i<other.size(); i++)
    {
        if((tempThis)[i] >= other[i])
        {
            result.push_back((tempThis)[i] - other[i]);
        }
        else
        {
            // 借位
            unsigned int j=0;
            for(j=i+1; j<tempThis.size(); j++)
            {
                if((tempThis)[j] > 0)
                {
                    break;
                }
            }
            (tempThis)[j] -= 1;
            for(unsigned k=j-1; k>i; k--)
            {
                (tempThis)[k] += 1;
            }
            if(j != i)
            {
                (tempThis)[i] = 2;
                (tempThis)[i] -= 1;
            }
            assert((tempThis)[i+1] >= 0);
            result.push_back(1);
        }
    }
}

```

```

    }
    if(i < tempThis.size())
    {
        for(; i < tempThis.size(); i++)
        {
            result.push_back((tempThis[i]));
        }
    }
    result.RemoveZero();
    return result;
}

friend ostream& operator << (ostream& o, const BitArray& a)
{
    for( int i = (int)a.size() - 1; i >= 0; i--)
    {
        o << (int)a[i];
    }
    o << endl;
    return o;
}

private:
void RemoveZero()
{
    unsigned int length = (unsigned int)this->size();
    for(int i = length - 1; i >= 0; i--)
    {
        if((*this)[i] == 0)
        {
            this->erase(this->begin() + i);
        }
        else
        {
            return;
        }
    }
}

void Smooth()
{
    for(unsigned int i = 0; i < this->size(); i++)
    {
        unsigned char a = (*this)[i];
        (*this)[i] = a % 2;
        if(i != this->size() - 1)
        {
            (*this)[i + 1] = (*this)[i + 1] + a / 2;
        }
        else
        {
            if(a / 2 == 0)
            {
                return;
            }
            else if(a / 2 == 1)
            {
                this->push_back(a / 2);
            }
            else
            {
                assert(false);
            }
        }
    }
}

```

```

    }
    }
}

};
//
// 求最接近 k 的 2 的 n 次方
//
int GetPowerOfTwo(int k)
{
    int result = 1;
    while( k > result)
    {
        result = result << 1;
    }
    return k;
}

void RECIPROCAL(BitArray& P, BitArray& Value)
{
    int k = (int) P.size();
    if(P.size() == 1)
    {
        // 返回 [1, 0]
        Value.clear();
        Value.resize(2);
        Value[0] = 0;
        Value[1] = 1;
    }
    else
    {
        int needLength = GetPowerOfTwo(k);
        // 如果 p 不是 2 的 n 次方, 那么插入额外的 0
        for(int i = 0; i < needLength - k; i++)
        {
            P.insert(P.begin(), 0);
        }
        int actualLength = k;
        k = needLength;
        int kDivideByTwo = k / 2;
        BitArray PDivideByTwo;
        PDivideByTwo.resize(kDivideByTwo);
        for(int i = (int)P.size() - 1; i >= (int)P.size() - kDivideByTwo; i--)
        {
            PDivideByTwo[i - (P.size() - kDivideByTwo)] = P[i];
        }
        BitArray C;
        RECIPROCAL(PDivideByTwo, C);

        // 加 (3 * k / 2) 个 0 到 C 的末尾, 得到  $C * 2^{(3 * k / 2)}$ 
        BitArray CTemp = C;
        for(int i = 0; i < 3 * k / 2; i++)
        {
            CTemp.insert(CTemp.begin(), 0);
        }

        BitArray D;
        D = CTemp - C * C * P;
        if(D.size() != 2 * k)
    }
}

```

```

    {
        // 补0
        D.push_back(0);
    }

    BitArray A;
    A.resize(k+1);
    for(int i=k; i>=0; i--)
    {
        A[i]=D[i+D.size()-(k+1)];
    }

    for(int i=2; i>=0; i--)
    {
        BitArray TwoPowerI;
        TwoPowerI.resize(i+1);
        TwoPowerI[i]=1;
        BitArray TwoPowerDoubleKMinusOne;
        TwoPowerDoubleKMinusOne.resize(2 * k);
        TwoPowerDoubleKMinusOne[2 * k-1]=1;

        if((A+TwoPowerI) * P <= TwoPowerDoubleKMinusOne)
        {
            A=A+TwoPowerI;
        }
    }
    Value=A;
    for(int i=0; i<needLength-actualLength; i++)
    {
        Value.erase(Value.begin());
    }
}

void TestRECIPROCAL()
{
    BitArray P;
    BitArray Value;

    P.resize(8);
    P[0]=1;
    P[3]=1;
    P[4]=1;
    P[7]=1;
    cout << "P" << endl;
    cout << P;
    RECIPROCAL(P, Value);
    cout << "Value" << endl;
    cout << Value;
}

int main()
{
    TestRECIPROCAL();
    system("pause");
    return 0;
}

```

A. 31 求多项式的倒数算法

```
#include <iostream>
```

```

#include <vector>
#include <assert.h>
using namespace std;
// RECIPROCAL
// 功能: 求多项式的倒数算法实现
// 函数: CoefficientArray RECIPROCAL(CoefficientArray& A)
// 输入:
//       CoefficientArray& A: 用多项式系数表示的多项式
// 输出:
//       CoefficientArray: 用多项式系数表示的多项式 A 的倒数
// 测试函数: TestRECIPROCAL()
// 输入:
//       A:  $X^7 - X^6 + X^5 + 2X^4 - X^3 - 3X^2 + X + 4$ 
//       表示为: [4, 1, -3, -1, 2, 1, -1, 1]
// 输出:
//        $X^7 + X^6 - 3X^4 - 4X^3 + 3X^2 + 15X + 12$ 
//       表示为: [12, 15, 3, -4, -3, 0, 1, 1]
// 使用多项式的系数来表示一个多项式, 高位在后, 低位在前
// 例如:  $X^7 - X^6 + X^5 + 2X^4 - X^3 - 3X^2 + X + 4$ 
// 表示为: [4, 1, -3, -1, 2, 1, -1, 1]
class CoefficientArray : public vector<double>
{
public:
    CoefficientArray& operator = (const CoefficientArray& other)
    {
        this->clear();
        this->resize(other.size());
        for(unsigned int i=0; i<other.size(); i++)
        {
            (*this)[i] = other[i];
        }
        return *this;
    }
    CoefficientArray operator* (double coeff)
    {
        CoefficientArray result = *this;
        for(unsigned i=0; i<result.size(); i++)
        {
            result[i] = result[i] * coeff;
        }
        result.RemoveZero();
        return result;
    }
    CoefficientArray operator* (const CoefficientArray& other)
    {
        CoefficientArray result;
        for(unsigned i=0; i<other.size(); i++)
        {
            CoefficientArray temp;
            for(unsigned j=0; j<i; j++)
            {
                temp.insert(temp.begin(), 0);
            }
            for(unsigned j=0; j<this->size(); j++)
            {
                temp.push_back(other[i] * (*this)[j]);
            }
            result = result + temp;
        }
    }
}

```

```

    result.RemoveZero();
    return result;
}

CoefficientArray operator - (const CoefficientArray& other)
{
    CoefficientArray result;
    CoefficientArray tempThis = * this;
    unsigned int i = 0;
    for(i=0; i<other.size() && i<this->size(); i++)
    {
        result.push_back((* this)[i] - other[i]);
    }
    if(i<this->size())
    {
        for(; i<this->size(); i++)
        {
            result.push_back((* this)[i]);
        }
    }
    else if(i<other.size())
    {
        for(; i<other.size(); i++)
        {
            result.push_back(-other[i]);
        }
    }
    result.RemoveZero();
    return result;
}

CoefficientArray operator + (const CoefficientArray& other)
{
    CoefficientArray a;
    unsigned i = 0;
    for (i=0; i<this->size() && i<other.size(); i++)
    {
        a.push_back((* this)[i] + other[i]);
    }
    if(i<this->size())
    {
        for(; i<this->size(); i++)
        {
            a.push_back((* this)[i]);
        }
    }
    else if(i<other.size())
    {
        for(; i<other.size(); i++)
        {
            a.push_back(other[i]);
        }
    }
    a.RemoveZero();
    return a;
}

friend ostream& operator << (ostream& o, const CoefficientArray& a)
{

```

```

    for( int i = a.size() - 1; i >= 0; i -- )
    {
        if(i == 1)
        {
            if(a[i] != 1)
            {
                o << a[i];
            }
            o << "x";
            if(a[0] >= 0)
            {
                o << " + ";
            }
            continue;
        }
        if(i == 0)
        {
            o << a[i];
            o << endl;
            continue;
        }
        if(a[i] != 1)
        {
            o << a[i];
        }
        if(i != 0)
        {
            o << "x" << i;
            if(a[i-1] >= 0)
            {
                o << " + ";
            }
        }
    }
    return o;
}

private:
void RemoveZero()
{
    unsigned int length = this->size();
    for(int i = length-1; i >= 0; i -- )
    {
        if((* this)[i] == 0)
        {
            this->erase(this->begin() + i);
        }
        else
        {
            return;
        }
    }
}

};

CoefficientArray RECIPROCAL(CoefficientArray& A)
{
    int k = (int)A.size();

    if(A.size() == 1)
    {
        CoefficientArray result;

```



```

        result.push_back(1.0 / A[0]);
        return result;
    }
    CoefficientArray AKDiviedByTwo;
    AKDiviedByTwo.resize(k / 2);
    for(unsigned i = k / 2; i < k; i++)
    {
        AKDiviedByTwo[i - k / 2] = A[i];
    }

    CoefficientArray Q = RECIPROCAL(AKDiviedByTwo);

    CoefficientArray Q1 = Q * 2;
    // 乘以  $X^{((3/2) * k - 2)}$ 
    // 由于k是2的n次方, 所有temp是个整数
    int temp = 3 * k / 2 - 2;
    for(int i = 0; i < temp; i++)
    {
        Q1.insert(Q1.begin(), 0);
    }
    CoefficientArray R = Q1 - Q * Q * A;
    if(k - 2 != 0)
    {
        R.erase(R.begin(), R.begin() + (k - 2));
    }
    return R;
}

void TestRECIPROCAL()
{
    CoefficientArray a;
    a.push_back(4);
    a.push_back(1);
    a.push_back(-3);
    a.push_back(-1);
    a.push_back(2);
    a.push_back(1);
    a.push_back(-1);
    a.push_back(1);
    cout << a;
    CoefficientArray S = RECIPROCAL(a);
    cout << S;
}

int main()
{
    TestRECIPROCAL();
    system("pause");
    return 0;
}

```

A.32 求一个整数的多个余数的算法

```

#include <iostream>
#include <vector>
#include <cmath>
using namespace std;
// RESIDUES
// 功能: 求一个整数的多个余数
// 函数: void RESIDUES(vector<int> &P, int u, vector<int> &results)
// 输入:

```

```

//          vector<int>& P: 一个素数的数组
//          int u: 整数 u
//  输出:
//          vector<int>& results: u 除以 P 中素数的余数
//  测试函数: TestRESIDUES()
//  输入:
//          P: [7, 5, 3, 2]
//          u: 201
//  输出:
//          results: [5, 1, 0, 1]
/*****
    矩阵类:
    封装一部分矩阵基本操作
*****/
class Matrix
{
public:
    Matrix(int row, int col)
    {
        int i=0;
        this->m_ppData=NULL;
        this->m_iRow=0;
        this->m_iCol=0;
        if(row<=0 || col<=0)
        {
            return;
        }
        this->m_iRow=row;
        this->m_iCol=col;
        this->m_ppData=new int* [row];
        for(i=0; i<row; i++)
        {
            this->m_ppData[i]=new int[col];
            memset(this->m_ppData[i], 0, sizeof(int) * this->m_iCol);
        }
    }
    Matrix(const Matrix& other)
    {
        this->operator=(other);
    }
    ~Matrix()
    {
        ClearMemory();
    }
    int* operator [] (int rowIndex)
    {
        int* pData=NULL;
        if(rowIndex>=0 && rowIndex<this->m_iRow)
        {
            if(this->m_ppData!=NULL)
            {
                pData=this->m_ppData[rowIndex];
            }
        }
        return pData;
    }
    int GetRowCount() { return m_iRow;}
    int GetColumnCount() { return m_iCol;}
private:
    void ClearMemory()

```

```

{
    int i=0;
    if(this->m_ppData != NULL)
    {
        for(i=0; i<this->m_iRow; i++)
        {
            delete [] (this->m_ppData[i]);
            this->m_ppData[i] = NULL;
        }
        delete [] this->m_ppData;
        this->m_ppData = NULL;
    }
}

private:
    int* * m_ppData;
    int m_iRow;
    int m_iCol;
};

//
// 计算以2为底的log
//
double log2(double a)
{
    return log(a) / log(2.0);
}

void RESIDUES(vector<int> &P, int u, vector<int> &results)
{
    // P的大小必须是2的n次方
    int k = (int)P.size();
    int t = (int)log2(k);
    Matrix Q(k, t);
    Matrix U(k, t+1);

    for(int i=0; i<k; i++)
    {
        Q[i][0] = P[i];
    }
    for(int j=1; j<t; j++)
    {
        for(int i=0; i<k; i += (int)pow(2.0, j))
        {
            Q[i][j] = Q[i][j-1] * Q[i + (int)pow(2.0, j-1)][j-1];
        }
    }

    U[0][t] = u;

    for(int j=t; j>0; j--)
    {
        for(int i=0; i<k; i += (int)pow(2.0, j))
        {
            U[i][j-1] = U[i][j] % Q[i][j-1];

            U[i + (int)pow(2.0, j-1)][j-1] = U[i][j] % Q[i + (int)pow(2.0, j-1)][j-1];
        }
    }
    for(int i=0; i<k; i++)
    {
        results.push_back((int)U[i][0]);
    }
}

```

```

    }
    return;
}
void TestRESIDUES()
{
    vector<int> P;
    vector<int> results;
    P.push_back(7);
    P.push_back(5);
    P.push_back(3);
    P.push_back(2);
    int u=201;
    RESIDUES(P, u, results);
    for(vector<int>:: iterator it = results.begin();
        it != results.end();
        it++)
    {
        cout << *it << endl;
    }
}
int main()
{
    TestRESIDUES();
    system("pause");
    return 0;
}

```

A.33 从一系列余数求原始数的算法

```

#include <iostream>
#include <vector>
#include <cmath>
using namespace std;
// ChineseRemainder
// 功能: 从一系列余数中求原始值
// 函数: int ChineseRemainder(vector<int> &P, vector<int> &D, vector<int> &U)
// 输入:
//       vector<int> &P: 一个素数的数组, P的大小是2的n次方
//       vector<int> &D: 一个整数数组,  $D_i = ((P/P_i)^{-1}) \% P_i$ 
//       vector<int> &U: u 除以P中素数的余数数组
// 输出:
//       int: 原始值
// 测试函数: TestChineseRemainder()
// 测试数据:
//       P: [2, 3, 5, 7]
//       U: [1, 2, 4, 3]
//       D: [1, 1, 3, 4]
// 测试结果:
//       原始值: 58
// *****
// 矩阵类:
// 封装一部分矩阵基本操作
// ***** /
class Matrix
{
public:
    Matrix(int row, int col)
    {
        int i=0;
        this->m_ppData = NULL;
    }

```

```

        this->m_iRow=0;
        this->m_iCol=0;
        if(row<=0 || col<=0)
        {
            return;
        }
        this->m_iRow=row;
        this->m_iCol=col;
        this->m_ppData=new int* [row];
        for(i=0; i<row; i++)
        {
            this->m_ppData[i]=new int[col];
            memset(this->m_ppData[i], 0, sizeof(int) * this->m_iCol);
        }
    }
    ~Matrix()
    {
        ClearMemory();
    }
    int* operator [] (int rowIndex)
    {
        int* pData=NULL;
        if(rowIndex>=0 && rowIndex<this->m_iRow)
        {
            if(this->m_ppData!=NULL)
            {
                pData=this->m_ppData[rowIndex];
            }
        }
        return pData;
    }
    int GetRowCount() { return m_iRow;}
    int GetColumnCount() { return m_iCol;}
private:
    void ClearMemory()
    {
        int i=0;
        if(this->m_ppData!=NULL)
        {
            for(i=0; i<this->m_iRow; i++)
            {
                delete [] (this->m_ppData[i]);
                this->m_ppData[i]=NULL;
            }
            delete [] this->m_ppData;
            this->m_ppData=NULL;
        }
    }
private:
    int* * m_ppData;
    int m_iRow;
    int m_iCol;
};
// 求D值
int GetD(int C, int P)
{
    int k=1;
    while(true)
    {
        if((k * C) % P == 1)

```

```

        {
            return k;
        }
        else
        {
            k++;
        }
    }
}
//
// 计算以 2 为底的 log
//
double log2(double a)
{
    return log(a) / log(2.0);
}
int ChineseRemainder(vector<int> &P, vector<int> &D, vector<int> &U)
{
    int k = (int)P.size();
    int t = (int)log2(k);
    Matrix S(k, t+1);
    Matrix Q(k, t+1);
    for(int i=0; i<k; i++)
    {
        Q[i][0] = P[i];
    }
    for(int j=1; j<=t; j++)
    {
        for(int i=0; i<k; i += (int)pow(2.0, j))
        {
            Q[i][j] = Q[i][j-1] * Q[i+(int)pow(2.0, j-1)][j-1];
        }
    }
    for(int i=0; i<k; i++)
    {
        S[i][0] = D[i] * U[i];
    }
    for(int j=1; j<=t; j++)
    {
        for(int i=0; i<k; i += (int)pow(2.0, j))
        {
            S[i][j] = S[i][j-1] * Q[i+(int)pow(2.0, j-1)][j-1] + S[i+(int)pow(2.0, j-1)]
[j-1] * Q[i][j-1];
        }
    }
    return S[0][t] % Q[0][t];
}
void TestChineseRemainder()
{
    vector<int> P;
    P.push_back(2);
    P.push_back(3);
    P.push_back(5);
    P.push_back(7);
    vector<int> U;
    U.push_back(1);
    U.push_back(2);
    U.push_back(4);
    U.push_back(3);
}

```

```

vector<int> D;
int c=1;
for(int i=0; i<(int)P.size(); i++)
{
    c=c * P[i];
}
for(int i=0; i<(int)P.size(); i++)
{
    D.push_back(GetD(c / P[i], P[i]));
}
int u=ChineseRemainder(P, D, U);
cout << u << endl;
}
int main()
{
    TestChineseRemainder();
    system("pause");
    return 0;
}

```

A.34 扩展的欧几里得算法

```

#include <iostream>
#include <vector>
#include <cmath>
using namespace std;
// ExtGCD
// 功能: 扩展的欧几里得算法
// 函数: void ExtGCD(int a0, int a1, int &a, int &x, int&y)
// 输入:
//      int a0: 正整数 a0
//      int a1: 正整数 a1
// 输出:
//      int a: a0 和 a1 的最大公约数
//      int x, y 整数 x 和 y, 使得  $a_0 * x + a_1 * y = a$ ;
// 测试函数: TestExtGCD()
// 输入:
//      a0: 57
//      a1: 33
// 输出:
//      a: 3
//      x: -4
//      y: 7
void ExtGCD(int a0, int a1, int &a, int &x, int&y)
{
    vector<int> A;
    vector<int> X;
    vector<int> Y;
    A.push_back(a0);
    A.push_back(a1);
    X.push_back(1);
    X.push_back(0);
    Y.push_back(0);
    Y.push_back(1);
    int i=1;
    while(A[i-1] % A[i] != 0)
    {
        int q=(int)floor((A[i-1] * 1.0 / A[i]));
        A.push_back(A[i-1]-q * A[i]);
        X.push_back(X[i-1]-q * X[i]);
    }
}

```

```

        Y.push_back(Y[i-1]-q * Y[i]);
        i=i+1;
    }
    a=A[i];
    x=X[i];
    y=Y[i];
    return;
}
void TestExtGCD()
{
    int a0=57;
    int a1=33;
    int a=0;
    int x=0;
    int y=0;
    ExtGCD(a0, a1, a, x, y);
    cout << a << endl;
    cout << x << endl;
    cout << y << endl;
}
int main()
{
    TestExtGCD();
    system("pause");
    return 0;
}

```

A.35 HGCD 算法

```

#include <iostream>
#include <vector>
#include <cmath>
#include <assert.h>
using namespace std;
// HGCD
// 功能: 求 2 个多项式的半公约数
// 函数: PolynomialMatrix HGCD(CoefficientArray& a0, CoefficientArray& a1)
// 输入:
//      CoefficientArray& a0: 用系数数组表示的多项式 a0
//      CoefficientArray& a1: 用系数数组表示的多项式 a1
// 输出:
//      PolynomialMatrix: 一个多项式的矩阵
// 测试函数: TestHGCD()
// 输入:
//      a0:  $X^5 + X^4 + X^3 + X^2 + X + 1$ 
//      a1:  $X^4 - 2X^3 + 3X^2 - X - 7$ 
// 输出:
//      |          1                      - (x+3)      |
//      |                                     |
//      |                                     |
//      | - (1/4 * x - 1/16)    1/4 * x^2 + 11/16 * x + 13/16 |
// 使用多项式的系数来表示一个多项式, 高位在后, 低位在前
// 例如:  $X^7 - X^6 + X^5 + 2X^4 - X^3 - 3X^2 + X + 4$ 
// 表示为: [4, 1, -3, -1, 2, 1, -1, 1]
//
class CoefficientArray : public vector <double>
{
public:
    int DEG() const
    {

```



```

    for(int i = (int)this->size() - 1; i >= 0; i++)
    {
        if((* this)[i] != 0)
        {
            return i;
        }
    }
    return 0;
}

CoefficientArray& operator = (const CoefficientArray& other)
{
    this->clear();
    this->resize(other.size());
    for(unsigned int i=0; i<other.size(); i++)
    {
        (* this)[i] = other[i];
    }
    return * this;
}

CoefficientArray operator* (double coeff)
{
    CoefficientArray result = * this;
    for(unsigned i=0; i<result.size(); i++)
    {
        result[i] = result[i] * coeff;
    }
    result.RemoveZero();
    return result;
}

CoefficientArray operator* (const CoefficientArray& other)
{
    CoefficientArray result;
    for(unsigned i=0; i<other.size(); i++)
    {
        CoefficientArray temp;
        for(unsigned j=0; j<i; j++)
        {
            temp.insert(temp.begin(), 0);
        }
        for(unsigned j=0; j<this->size(); j++)
        {
            temp.push_back(other[i] * (* this)[j]);
        }

        result = result + temp;
    }
    result.RemoveZero();
    return result;
}

CoefficientArray operator % (const CoefficientArray& other)
{
    CoefficientArray result;
    this->DivideBy(other, result);
    return result;
}

CoefficientArray operator /(const CoefficientArray& other)
{
    CoefficientArray result;
    CoefficientArray residue;
    result = this->DivideBy(other, residue);
}

```

```

    return result;
}

bool CanDivideBy(CoefficientArray& other)
{
    CoefficientArray moduloResult = * this % other;
    if(moduloResult.size() == 0)
    {
        return true;
    }
    return false;
}

CoefficientArray operator - (const CoefficientArray& other)
{
    CoefficientArray result;
    CoefficientArray tempThis = * this;
    unsigned int i = 0;
    for(i=0; i<other.size() && i<this->size(); i++)
    {
        result.push_back((* this)[i] - other[i]);
    }
    if(i<this->size())
    {
        for(; i<this->size(); i++)
        {
            result.push_back((* this)[i]);
        }
    }
    else if(i<other.size())
    {
        for(; i<other.size(); i++)
        {
            result.push_back(-other[i]);
        }
    }
    result.RemoveZero();
    return result;
}

CoefficientArray operator + (const CoefficientArray& other)
{
    CoefficientArray a;
    unsigned i = 0;
    for(i=0; i<this->size() && i<other.size(); i++)
    {
        a.push_back((* this)[i] + other[i]);
    }
    if(i<this->size())
    {
        for(; i<this->size(); i++)
        {
            a.push_back((* this)[i]);
        }
    }
    else if(i<other.size())
    {
        for(; i<other.size(); i++)
        {
            a.push_back(other[i]);
        }
    }
}

```

```

    }
}
a.RemoveZero();
return a;
}

friend ostream& operator << (ostream& o, CoefficientArray& a)
{
    for( int i = (int)a.size() - 1; i >= 0; i -- )
    {
        if(i != 1)
        {
            if(a[i] != 1)
            {
                o << a[i];
            }
            o << "x";
            if(a[0] >= 0)
            {
                o << " + ";
            }
            continue;
        }
        if(i == 0)
        {
            o << a[i];
            o << endl;
            continue;
        }
        if(a[i] != 1)
        {
            o << a[i];
        }
        if( i != 0)
        {
            o << "x" << i;
            if(a[i - 1] >= 0)
            {
                o << " + ";
            }
        }
    }
    return o;
}

private:

CoefficientArray DivideBy(const CoefficientArray& divisor, CoefficientArray& residue)
{
    CoefficientArray result;
    CoefficientArray dividend = * this;
    residue = * this;
    if(this->size() == 0 || this->DEG() < divisor.DEG())
    {
        return result;
    }
    do
    {
        dividend = residue;
        int m = dividend.DEG();
        int n = divisor.DEG();
    }
}

```

```

        CoefficientArray quotient;
        quotient.resize(m-n+1);
        quotient[m-n] = dividend[m] / divisor[n];
        residue = dividend - quotient * divisor;
        result = result + quotient;
    }while(! (residue.size() == 0 || residue.DEG() < divisor.DEG()));
    return result;
}

void RemoveZero()
{
    int length = (int)this->size();
    for(int i=length-1; i >= 0; i--)
    {
        if((*this)[i] == 0)
        {
            this->erase(this->begin() + i);
        }
        else
        {
            return;
        }
    }
}

};

//
// 多项式矩阵
//
class PolynomialMatrix
{
public:
    PolynomialMatrix(int row, int col)
    {
        int i=0;
        this->m_ppData = NULL;
        this->m_iRow = 0;
        this->m_iCol = 0;
        if(row <= 0 || col <= 0)
        {
            return;
        }
        this->m_iRow = row;
        this->m_iCol = col;
        this->m_ppData = new CoefficientArray* [row];
        for(i = 0; i < row; i++)
        {
            this->m_ppData[i] = new CoefficientArray[col];
        }
    }
    PolynomialMatrix(const PolynomialMatrix& other)
    {
        this->operator = (other);
    }
    ~PolynomialMatrix()
    {
        ClearMemory();
    }
    CoefficientArray* operator [] (int rowIndex)
    {

```

```

CoefficientArray* pData = NULL;
if (rowIndex >= 0 && rowIndex < this->m_iRow)
{
    if (this->m_ppData != NULL)
    {
        pData = this->m_ppData[rowIndex];
    }
}
return pData;
}

PolynomialMatrix operator * (const PolynomialMatrix& other)
{
    int i=0;
    int j=0;
    int k=0;
    PolynomialMatrix result(this->m_iRow, other.m_iCol);
    if (this->m_iCol != other.m_iRow)
    {
        assert(0);
        return result;
    }
    for (i=0; i<result.m_iRow; i++)
    {
        for (j=0; j<result.m_iCol; j++)
        {
            for (k=0; k<this->m_iCol; k++)
            {
                result[i][j] = result[i][j] +
                    const_cast<PolynomialMatrix&>(*this)[i][k] *
                    const_cast<PolynomialMatrix&>(other)[k][j];
            }
        }
    }
    return result;
}

PolynomialMatrix& operator = (const PolynomialMatrix& other)
{
    if (&other == this)
    {
        return *this;
    }
    int i=0;
    this->m_iCol = other.m_iCol;
    this->m_iRow = other.m_iRow;
    this->m_ppData = new CoefficientArray* [this->m_iRow];
    for (i=0; i<this->m_iRow; i++)
    {
        this->m_ppData[i] = new CoefficientArray[this->m_iCol];
        for (int j=0; j<this->m_iCol; j++)
        {
            (const_cast<PolynomialMatrix*>(this)->m_ppData[i])[j] =
                const_cast<PolynomialMatrix&>(other)[i][j];
        }
    }
    return *this;
}

friend ostream& operator << (ostream& o, const PolynomialMatrix& a)
{
    for (int i=0; i<a.m_iRow; i++)
    {

```

```

        for(int j=0; j<a.m_iCol; j++)
        {
            o << a.m_ppData[i][j] << endl;
        }
    }
    return o;
}
int GetRowCount() { return m_iRow;}
int GetColumnCount() { return m_iCol;}
private:
void ClearMemory()
{
    int i=0;
    if(this->m_ppData != NULL)
    {
        for(i=0; i<this->m_iRow; i++)
        {
            delete [] (this->m_ppData[i]);
            this->m_ppData[i] = NULL;
        }
        delete [] this->m_ppData;
        this->m_ppData = NULL;
    }
}
private:
CoefficientArray* * m_ppData;
int m_iRow;
int m_iCol;
};
PolynomialMatrix HGCD(CoefficientArray& a0, CoefficientArray& a1)
{
    int m = (int) floor(a0.DEG() / 2.0);
    if(a1.DEG() < = m)
    {
        PolynomialMatrix result(2, 2);
        result[0][0].push_back(1);
        result[1][1].push_back(1);
        return result;
    }
    else
    {
        CoefficientArray b0;
        CoefficientArray c0;
        for(int i=0; i<m; i++)
        {
            c0.push_back(a0[i]);
        }
        for(int i=m; i<(int)a0.size(); i++)
        {
            b0.push_back(a0[i]);
        }
        CoefficientArray b1;
        CoefficientArray c1;
        for(int i=0; i<m; i++)
        {
            c1.push_back(a1[i]);
        }
        for(int i=m; i<(int)a1.size(); i++)
        {

```

```

        bl.push_back(a1[i]);
    }
    PolynomialMatrix R = HGCD(b0, b1);
    PolynomialMatrix MatrixA0AndA1(2, 1);
    MatrixA0AndA1[0][0] = a0;
    MatrixA0AndA1[1][0] = a1;
    PolynomialMatrix MatrixDAndE = R * MatrixA0AndA1;
    CoefficientArray d = MatrixDAndE[0][0];
    CoefficientArray e = MatrixDAndE[1][0];
    CoefficientArray f = d % e;
    int mDivideTwo = m / 2;
    CoefficientArray g0;
    CoefficientArray h0;
    for(int i = 0; i < mDivideTwo; i++)
    {
        h0.push_back(e[i]);
    }

    for(int i = mDivideTwo; i < (int)e.size(); i++)
    {
        g0.push_back(e[i]);
    }
    CoefficientArray g1;
    CoefficientArray h1;
    for(int i = 0; i < mDivideTwo; i++)
    {
        h1.push_back(f[i]);
    }
    for(int i = mDivideTwo; i < (int)f.size(); i++)
    {
        g1.push_back(f[i]);
    }
    PolynomialMatrix S = HGCD(g0, g1);
    CoefficientArray q = d / e;
    PolynomialMatrix M(2, 2);
    M[0][1].push_back(1);
    M[1][0].push_back(1);
    for(int i = 0; i < (int)q.size(); i++)
    {
        q[i] = -q[i];
    }
    M[1][1] = q;
    return S * M * R;
}

void TestHGCD()
{
    CoefficientArray p1;
    p1.push_back(1);
    p1.push_back(1);
    p1.push_back(1);
    p1.push_back(1);
    p1.push_back(1);
    p1.push_back(1);
    CoefficientArray p2;
    p2.push_back(-7);
    p2.push_back(-1);
    p2.push_back(3);
    p2.push_back(-2);
    p2.push_back(1);
}

```

```

    PolynomialMatrix result = HGCD(p1, p2);
    cout << result << endl;
}
int main()
{
    TestHGCD();
    system("pause");
    return 0;
}

```

A.36 求两个多项式的公约数

```

#include <iostream>
#include <vector>
#include <cmath>
using namespace std;
// GCD
// 功能: 求两个多项式的公约数
// 函数: CoefficientArray GCD(CoefficientArray& a0, CoefficientArray& a1)
// 输入:
//      CoefficientArray& a0: 用系数数组表示的多项式 a0
//      CoefficientArray& a1: 用系数数组表示的多项式 a1
// 输出:
//      CoefficientArray: a0 和 a1 的多项式
// 测试函数: TestGCD()
// 输入:
//      a0:  $X^5 + X^4 + X^3 + X^2 + X + 1$ 
//      a1:  $X^4 - 2X^3 + 3X^2 - X - 7$ 
// 输出:
//       $3952x + 3952$ 
// 使用多项式的系数来表示一个多项式, 高位在后, 低位在前
// 例如:  $X^7 - X^6 + X^5 + 2X^4 - X^3 - 3X^2 + X + 4$ 
// 表示为: [4, 1, -3, -1, 2, 1, -1, 1]
//
class CoefficientArray : public vector<double>
{
public:
    int DEG() const
    {
        for(int i = (int)this->size() - 1; i >= 0; i++)
        {
            if((*this)[i] != 0)
            {
                return i;
            }
        }
        return 0;
    }
    CoefficientArray& operator = (const CoefficientArray& other)
    {
        this->clear();
        this->resize(other.size());
        for(unsigned int i = 0; i < other.size(); i++)
        {
            (*this)[i] = other[i];
        }
        return *this;
    }
    CoefficientArray operator* (double coeff)
    {

```



```

CoefficientArray result = * this;
for(unsigned i=0; i<result.size(); i++)
{
    result[i] = result[i] * coeff;
}
result.RemoveZero();
return result;
}

CoefficientArray operator* (const CoefficientArray& other)
{
    CoefficientArray result;
    for(unsigned i=0; i<other.size(); i++)
    {
        CoefficientArray temp;
        for(unsigned j=0; j<i; j++)
        {
            temp.insert(temp.begin(), 0);
        }
        for(unsigned j=0; j<this->size(); j++)
        {
            temp.push_back(other[i] * (*this)[j]);
        }

        result = result + temp;
    }
    result.RemoveZero();
    return result;
}

CoefficientArray operator % (const CoefficientArray& other)
{
    CoefficientArray result;
    this->DivideBy(other, result);
    return result;
}

CoefficientArray operator / (const CoefficientArray& other)
{
    CoefficientArray result;
    CoefficientArray residue;
    result = this->DivideBy(other, residue);
    return result;
}

bool CanDivideBy(CoefficientArray& other)
{
    CoefficientArray moduloResult = * this % other;
    if(moduloResult.size() == 0)
    {
        return true;
    }
    return false;
}

CoefficientArray operator - (const CoefficientArray& other)
{
    CoefficientArray result;
    CoefficientArray tempThis = * this;
    unsigned int i = 0;
    for(i=0; i<other.size() && i<this->size(); i++)
    {
        result.push_back((*this)[i] - other[i]);
    }

```

```

    }
    if(i < this->size())
    {
        for(; i < this->size(); i++)
        {
            result.push_back((*this)[i]);
        }
    }
    else if(i < other.size())
    {
        for(; i < other.size(); i++)
        {
            result.push_back(-other[i]);
        }
    }
    result.RemoveZero();
    return result;
}

CoefficientArray operator + (const CoefficientArray& other)
{
    CoefficientArray a;
    unsigned i = 0;
    for( i = 0; i < this->size() && i < other.size(); i++)
    {
        a.push_back((*this)[i] + other[i]);
    }
    if(i < this->size())
    {
        for(; i < this->size(); i++)
        {
            a.push_back((*this)[i]);
        }
    }
    else if(i < other.size())
    {
        for(; i < other.size(); i++)
        {
            a.push_back(other[i]);
        }
    }
    a.RemoveZero();
    return a;
}

friend ostream& operator << (ostream& o, CoefficientArray& a)
{
    for( int i = (int)a.size() - 1; i >= 0; i--)
    {
        if(i == 1)
        {
            if(a[i] != 1)
            {
                o << a[i];
            }
            o << "x";
            if(a[0] >= 0)
            {
                o << "+";
            }
        }
    }
}

```

```

        continue;
    }
    if(i == 0)
    {
        o << a[i];
        o << endl;
        continue;
    }
    if(a[i] != 1)
    {
        o << a[i];
    }
    if(i != 0)
    {
        o << "x" << i;
        if(a[i-1] >= 0)
        {
            o << " + ";
        }
    }
    }
    return o;
}
private:

CoefficientArray DivideBy(const CoefficientArray& divisor, CoefficientArray& residue)
{
    CoefficientArray result;
    CoefficientArray dividend = * this;
    residue = * this;
    if(this->size() == 0 || this->DEG() < divisor.DEG())
    {
        return result;
    }
    do
    {
        dividend = residue;
        int m = dividend.DEG();
        int n = divisor.DEG();
        CoefficientArray quotient;
        quotient.resize(m - n + 1);
        quotient[m - n] = dividend[m] / divisor[n];
        residue = dividend - quotient * divisor;
        result = result + quotient;
    } while(! (residue.size() == 0 || residue.DEG() < divisor.DEG()));
    return result;
}

void RemoveZero()
{
    int length = (int) this->size();
    for(int i = length - 1; i >= 0; i--)
    {
        if((* this)[i] == 0)
        {
            this->erase(this->begin() + i);
        }
        else
        {
            return;
        }
    }
}

```

```

    }
}
};

//
// 多项式矩阵
//
class PolynomialMatrix
{
public:
    PolynomialMatrix(int row, int col)
    {
        int i=0;
        this->m_ppData=NULL;
        this->m_iRow=0;
        this->m_iCol=0;
        if(row<=0 || col<=0)
        {
            return;
        }
        this->m_iRow=row;
        this->m_iCol=col;
        this->m_ppData=new CoefficientArray* [row];
        for(i=0; i<row; i++)
        {
            this->m_ppData[i]=new CoefficientArray[col];
        }
    }
    PolynomialMatrix(const PolynomialMatrix& other)
    {
        this->operator=(other);
    }
    ~PolynomialMatrix()
    {
        ClearMemory();
    }
    CoefficientArray* operator [] (int rowIndex) const
    {
        CoefficientArray* pData=NULL;
        if(rowIndex>=0 && rowIndex<this->m_iRow)
        {
            if(this->m_ppData!=NULL)
            {
                pData=this->m_ppData[rowIndex];
            }
        }
        return pData;
    }
    PolynomialMatrix operator * (PolynomialMatrix& other)
    {
        int i=0;
        int j=0;
        int k=0;
        PolynomialMatrix result(this->m_iRow, other.m_iCol);
        if(this->m_iCol!=other.m_iRow)
        {
            return result;
        }
        for(i=0; i<result.m_iRow; i++)

```

```

    {
        for(j=0; j<result.m_iCol; j++)
        {
            for(k=0; k<this->m_iCol; k++)
            {
                result[i][j] = result[i][j] + (*this)[i][k] * other[k][j];
            }
        }
    }
    return result;
}
PolynomialMatrix& operator = (const PolynomialMatrix& other)
{
    if(&other == this)
    {
        return *this;
    }
    int i=0;
    this->m_iCol = other.m_iCol;
    this->m_iRow = other.m_iRow;
    this->m_ppData = new CoefficientArray* [this->m_iRow];
    for(i=0; i<this->m_iRow; i++)
    {
        this->m_ppData[i] = new CoefficientArray[this->m_iCol];
        for(int j=0; j<this->m_iCol; j++)
        {
            (this->m_ppData[i])[j] = other[i][j];
        }
    }
    return *this;
}
int GetRowCount() { return m_iRow;}
int GetColumnCount() { return m_iCol;}
private:
void ClearMemory()
{
    int i=0;
    if(this->m_ppData != NULL)
    {
        for(i=0; i<this->m_iRow; i++)
        {
            delete [] (this->m_ppData[i]);
            this->m_ppData[i] = NULL;
        }
        delete [] this->m_ppData;
        this->m_ppData = NULL;
    }
}
private:
CoefficientArray* * m_ppData;
int m_iRow;
int m_iCol;
};
PolynomialMatrix HGCD(CoefficientArray& a0, CoefficientArray& a1)
{
    int m = (int) floor(a0.DEG() / 2.0);
    if(a1.DEG() < m)
    {
        PolynomialMatrix result(2, 2);
        result[0][0].push_back(1);
    }
}

```

```

    result[1][1].push_back(1);
    return result;
}
else
{
    CoefficientArray b0;
    CoefficientArray c0;
    for(int i=0; i<m; i++)
    {
        c0.push_back(a0[i]);
    }
    for(int i=m; i<(int)a0.size(); i++)
    {
        b0.push_back(a0[i]);
    }
    CoefficientArray b1;
    CoefficientArray c1;
    for(int i=0; i<m; i++)
    {
        c1.push_back(a1[i]);
    }
    for(int i=m; i<(int)a1.size(); i++)
    {
        b1.push_back(a1[i]);
    }
    PolynomialMatrix R=HGCD(b0, b1);
    PolynomialMatrix MatrixA0AndA1(2, 1);
    MatrixA0AndA1[0][0]=a0;
    MatrixA0AndA1[1][0]=a1;
    PolynomialMatrix MatrixDAndE=R * MatrixA0AndA1;
    CoefficientArray d=MatrixDAndE[0][0];
    CoefficientArray e=MatrixDAndE[1][0];
    CoefficientArray f=d % e;
    int mDivideTwo=m / 2;
    CoefficientArray g0;
    CoefficientArray h0;
    for(int i=0; i<mDivideTwo; i++)
    {
        h0.push_back(e[i]);
    }

    for(int i=mDivideTwo; i<(int)e.size(); i++)
    {
        g0.push_back(e[i]);
    }
    CoefficientArray g1;
    CoefficientArray h1;
    for(int i=0; i<mDivideTwo; i++)
    {
        h1.push_back(f[i]);
    }
    for(int i=mDivideTwo; i<(int)f.size(); i++)
    {
        g1.push_back(f[i]);
    }
    PolynomialMatrix S=HGCD(g0, g1);
    CoefficientArray q=d / e;
    PolynomialMatrix M(2, 2);
    M[0][1].push_back(1);

```

```

        M[1][0].push_back(1);
        for(int i=0; i < (int)q.size(); i++)
        {
            q[i] = -q[i];
        }
        M[1][1] = q;
        return S * M * R;
    }
}

CoefficientArray GCD(CoefficientArray& a0, CoefficientArray& a1)
{
    if(a0.CanDivideBy(a1))
    {
        return a1;
    }

    PolynomialMatrix R = HGCD(a0, a1);
    PolynomialMatrix a(2, 1);
    a[0][0] = a0;
    a[1][0] = a1;

    PolynomialMatrix b(2, 1);
    b = R * a;
    CoefficientArray b0 = b[0][0];
    CoefficientArray b1 = b[1][0];
    if(b0.CanDivideBy(b1))
    {
        return b1;
    }
    else
    {
        CoefficientArray c = b0 % b1;
        return GCD(b1, c);
    }
}

void TestGCD()
{
    CoefficientArray p1;
    p1.push_back(1);
    p1.push_back(1);
    p1.push_back(1);
    p1.push_back(1);
    p1.push_back(1);
    p1.push_back(1);
    CoefficientArray p2;
    p2.push_back(-7);
    p2.push_back(-1);
    p2.push_back(3);
    p2.push_back(-2);
    p2.push_back(1);
    CoefficientArray result = GCD(p1, p2);
    cout << "result:" << result << endl;
}

int main()
{
    TestGCD();
    system("pause");
    return 0;
}

```

A.37 NFA 的模拟

```

#include <iostream>
#include <vector>
#include <set>
#include <queue>
#include <string>
using namespace std;
// SimulationOfNFA
// 功能: NFA 的模拟
// 函数: void SimulationOfNFA(set <int> & S, vector <char> & I, StateTransition& transi-
tion, int s0, int F, string& x, vector <set <int> > & SSequence)
// 输入:
//      set <int> & S: NFA M 的所有状态
//      vector <char> & I: NFA M 的字母表
//      StateTransition& transition: NFA M 的转换函数
//      int s0: NFA M 的初始状态
//      int F: NFA M 的结束状态
// 输出:
//      vector <set <int> > & SSequence: 状态转换的序列
// 测试函数: TestSimulationOfNFA()
//      使用图 9.6 中的例子 ab* + c
// 输入:
//      S: [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]
//      I: [a, b, c]
//      transition: ...
//      s0: 1
//      F: 10
// 输出:
//      SSequence: {[1, 2, 8], [3, 4, 5, 7, 10], [5, 6, 7, 10]}
//      用来表示 NFA 的转换函数,
//      使用 '\0' 表示空输入
//
class StateTransition
{
private:
    class StateTransitionItem
    {
public:
        set <int> states;
        int state;
        char character;
    };
private:
    vector <StateTransitionItem* > items;
public:
    StateTransition()
    {
    }
    ~StateTransition()
    {
        for(int i=0; i < (int)this->items.size(); i++)
        {
            delete this->items[i];
        }
    }
public:
    void AddToTransition(int state, char character, set <int> states)
    {

```



```

    for(int i=0; i < (int)this->items.size(); i++)
    {
        StateTransitionItem* pItem=this->items[i];
        if(pItem->state == state && pItem->charater == charater)
        {
            pItem->states = states;
            return;
        }
    }
    StateTransitionItem* pItem=new StateTransitionItem();
    pItem->states = states;
    pItem->state = state;
    pItem->charater = charater;
    this->items.push_back(pItem);
}
set<int> operator()(int state, char charater)
{
    set<int> empty;
    for(int i=0; i < (int)this->items.size(); i++)
    {
        StateTransitionItem* pItem=this->items[i];
        if(pItem->state == state && pItem->charater == charater)
        {
            return pItem->states;
        }
    }
    return empty;
}
};

void SimulationOfNdfa(set<int> & S, vector<char> & I, StateTransition& transition, int s0,
int F, string& x, vector<set<int>> & SSequence)
{
    if(x.size() == 0)
    {
        return;
    }
    set<int>::iterator setIterator;
    SSequence.resize(x.size()+1);

    for(int i=0; i <= (int)x.size(); i++)
    {
        if(i == 0)
        {
            SSequence[0].insert(s0);
        }
        else
        {
            for(setIterator = SSequence[i-1].begin();
                setIterator != SSequence[i-1].end();
                setIterator++)
            {
                set<int> states = transition(* setIterator, x[i-1]);
                for(set<int>::iterator it = states.begin(); it != states.end(); it++)
                {
                    SSequence[i].insert(* it);
                }
            }
        }
    }
    vector<bool> consideredVec;

```

```

consideredVec.resize(S.size() + 1);
for(setIterator = SSequence[i].begin(); setIterator != SSequence[i].end(); setIter-
ator++)
{
    consideredVec[* setIterator] = true;
}

queue<int> QUEUE;
for(set<int>::iterator it = SSequence[i].begin(); it != SSequence[i].end(); it++)
{
    QUEUE.push(* it);
}
while(QUEUE.empty() == false)
{
    int t = QUEUE.front();
    QUEUE.pop();
    set<int> states = transition(t, '\0');
    for(set<int>::iterator it = states.begin(); it != states.end(); it++)
    {
        if(consideredVec[* it] == false)
        {
            consideredVec[* it] = true;
            QUEUE.push(* it);
            SSequence[i].insert(* it);
        }
    }
}
}

void TestSimulationOfNFA()
{
    set<int> S;
    S.insert(1);
    S.insert(2);
    S.insert(3);
    S.insert(4);
    S.insert(5);
    S.insert(6);
    S.insert(7);
    S.insert(8);
    S.insert(9);
    S.insert(10);
    vector<char> I;
    I.push_back('a');
    I.push_back('b');
    I.push_back('c');
    StateTransition transition;
    set<int> empty;
    set<int> nextStates;
    transition.AddToTransition(1, 'a', empty);
    transition.AddToTransition(1, 'b', empty);
    transition.AddToTransition(1, 'c', empty);
    nextStates.clear();
    nextStates.insert(2);
    nextStates.insert(8);
    transition.AddToTransition(1, '\0', nextStates);
    nextStates.clear();
    nextStates.insert(3);
    transition.AddToTransition(2, 'a', nextStates);
    transition.AddToTransition(2, 'b', empty);

```

```

transition.AddToTransition(2, 'c', empty);
transition.AddToTransition(2, '\0', empty);
transition.AddToTransition(3, 'a', empty);
transition.AddToTransition(3, 'b', empty);
transition.AddToTransition(3, 'c', empty);
nextStates.clear();
nextStates.insert(4);
transition.AddToTransition(3, '\0', nextStates);
transition.AddToTransition(4, 'a', empty);
transition.AddToTransition(4, 'b', empty);
transition.AddToTransition(4, 'c', empty);
nextStates.clear();
nextStates.insert(5);
nextStates.insert(7);
transition.AddToTransition(4, '\0', nextStates);
transition.AddToTransition(5, 'a', empty);
nextStates.clear();
nextStates.insert(6);
transition.AddToTransition(5, 'b', nextStates);
transition.AddToTransition(5, 'c', empty);
transition.AddToTransition(5, '\0', empty);
transition.AddToTransition(6, 'a', empty);
transition.AddToTransition(6, 'b', empty);
transition.AddToTransition(6, 'c', empty);
nextStates.clear();
nextStates.insert(5);
nextStates.insert(7);
transition.AddToTransition(6, '\0', nextStates);
transition.AddToTransition(7, 'a', empty);
transition.AddToTransition(7, 'b', empty);
transition.AddToTransition(7, 'c', empty);
nextStates.clear();
nextStates.insert(10);
transition.AddToTransition(7, '\0', nextStates);
transition.AddToTransition(8, 'a', empty);
transition.AddToTransition(8, 'b', empty);
nextStates.clear();
nextStates.insert(9);
transition.AddToTransition(8, 'c', nextStates);
transition.AddToTransition(8, '\0', empty);
transition.AddToTransition(9, 'a', empty);
transition.AddToTransition(9, 'b', empty);
transition.AddToTransition(9, 'c', empty);
nextStates.clear();
nextStates.insert(10);
transition.AddToTransition(9, '\0', nextStates);
transition.AddToTransition(10, 'a', empty);
transition.AddToTransition(10, 'b', empty);
transition.AddToTransition(10, 'c', empty);
transition.AddToTransition(10, '\0', empty);
int s0 = 1;
int F = 10;
string x = "ab";
vector<set<int>> SSequence;
SimulationOfNDFA(S, I, transition, s0, F, x, SSequence);
for(vector<set<int>>::iterator vecit = SSequence.begin();
    vecit != SSequence.end(); vecit++)
{
    cout << "[";
    for(set<int>::iterator setit = vecit->begin();

```

```

        setit != vecit -> end();
        setit ++ )
    {
        cout << * setit << " \t";
    }
    cout << "]" << endl;
}
}
int main()
{
    TestSimulationOfNdfa();
    system("pause");
    return 0;
}

```

A. 38 计算失败函数的算法

```

#include <iostream>
#include <string>
#include <vector>
using namespace std;
// FAILURE
// 功能: 求一个模式串的失败函数
// 函数: void FAILURE(string& patternB, vector<int> & f)
// 输入:
//       string& patternB: 模式串 B
// 输出:
//       vector<int> & f: B 的失败函数
// 测试函数: TestFAILURE()
// 输入:
//       patternB : "aabbaab";
// 输出:
//       f: [0, 1, 0, 0, 1, 2, 3]
//
// 数组的第一个元素没有使用
//
void FAILURE(const string& patternB, vector<int> & f)
{
    if(patternB.size() == 0)
    {
        return;
    }

    f.clear();
    f.resize(patternB.size());

    f[1] = 0;
    for(int j = 2; j < (int)patternB.size(); j++)
    {
        int i = f[j-1];
        while(patternB[j] != patternB[i+1] && i > 0)
        {
            i = f[i];
        }
        if(patternB[j] == patternB[i+1] && i == 0)
        {
            f[j] = 0;
        }
        else
    }
}

```

```

        {
            f[j] = i + 1;
        }
    }
}

void TestFAILURE()
{
    string pattern = "aabbaab";
    vector<int> f;
    FAILURE(pattern, f);
    for(vector<int>::iterator it = f.begin();
        it != f.end();
        it++)
    {
        cout << *it << endl;
    }
}

int main()
{
    TestFAILURE();
    system("pause");
    return 0;
}

```

A.39 判断一个输入串是否与模式匹配的算法

```

#include <iostream>
#include <vector>
#include <set>
#include <queue>
#include <string>
#include <map>
using namespace std;
// SimulationOf2DPDA
// 功能: 模拟 2DPDA, 判断一个输入串是否符合模式
// 函数: bool SimulationOf2DPDA(vector<Configuration*> &configurations)
// 输入:
// (vector<Configuration*> &configurations: 由输入串和 2DPDA 构成的所有配置
// 输出:
// bool: 输入串符合模式则返回 true, 否则返回 false
// 测试函数: TestFAILURE()
// 输入:
// 图 9.17 代表的配置
// 输出:
// true
// 备注:
// 为了简化程序, 这里省略了自动机和输入串生成所有配置的过程。
typedef struct Configuration
{
    string strConfigName;
    int iNumberOfMove;
}Configuration;
typedef map<Configuration*, Configuration*> TERM;
typedef map<Configuration*, set<Configuration*>> PRED;
typedef set<pair<Configuration*, Configuration*>> NEW;
typedef set<Configuration*> TEMP;
void UPDATE(Configuration* C, Configuration* D, TERM& Term, PRED& Pred, NEW& New, TEMP&
Temp);
bool SimulationOf2DPDA(vector<Configuration*> &configurations)
{

```

```

TERM Term;
PRED Pred;
NEW New;
TEMP Temp;
set <Configuration* > empty;
// 初始化
for(int i=0; i < (int)configurations.size(); i++)
{
    Configuration* C = configurations[i];
    Term[C] = NULL;
    Pred[C] = empty;
}
vector <Configuration* > terminalConfigs;
// 找到所有的结束配置
for(int i=0; i < (int)configurations.size(); i++)
{
    if(i != configurations.size() - 1)
    {
        if(configurations[i] -> iNumberOfMove > configurations[i + 1] -> iNumberOfMove)
        {
            terminalConfigs.push_back(configurations[i]);
        }
    }
    else
    {
        terminalConfigs.push_back(configurations[i]);
    }
}
for(int i=0; i < (int)terminalConfigs.size(); i++)
{
    Configuration* C = terminalConfigs[i];
    Term[C] = C;
    New.insert(make_pair(C, C));
}
// 找到所有的在一步内, move 数不变的配置对
for(int i=0; i < (int)configurations.size(); i++)
{
    if(i != configurations.size() - 1)
    {
        if(configurations[i] -> iNumberOfMove == configurations[i + 1] -> iNumberOfMove)
        {
            UPDATE(configurations[i], configurations[i + 1], Term, Pred, New, Temp);
        }
    }
}
while(! New.empty())
{
    pair<Configuration*, Configuration* > p = * (New.begin());
    New.erase(New.begin());

    Configuration* CComma = p.first;
    Configuration* DComma = p.second;
    int iCCommaIndex = -1;
    int iDCommaIndex = -1;
    for(int i=0; i < (int)configurations.size(); i++)
    {
        if(configurations[i] == CComma)
        {
            iCCommaIndex = i;
        }
    }
}

```

```

        if(configurations[i] == DComma)
        {
            iDCommaIndex = i;
        }
    }
    // 判断是否存在比(C, D)低的配置对
    if(
        iCCommaIndex != 0
        && iDCommaIndex != (configurations.size() - 1)
        && configurations[iCCommaIndex - 1] -> iNumberOfMove < configurations[iCCommaIndex - 1] -> iNumberOfMove
        && configurations[iDCommaIndex] -> iNumberOfMove > configurations[iDCommaIndex + 1] -> iNumberOfMove
    )
    {
        Configuration* C = configurations[iCCommaIndex - 1];
        Configuration* D = configurations[iDCommaIndex + 1];
        UPDATE(C, D, Term, Pred, New, Temp);
    }
    return Term[configurations[0]] == configurations[configurations.size() - 1];
}
void UPDATE(Configuration* C, Configuration* D, TERM& Term, PRED& Pred, NEW& New, TEMP& Temp)
{
    if (Term[D] == NULL)
    {
        set<Configuration*> configSet = Pred[D];
        configSet.insert(C);
        Pred[D] = configSet;
    }
    else
    {
        Temp.insert(C);
        while(Temp.empty() == false)
        {
            Configuration* B = * (Temp.begin());
            Temp.erase(Temp.begin());
            Term[B] = Term[D];
            New.insert(make_pair(B, Term[D]));
            set<Configuration*> predB = Pred[B];
            for(set<Configuration*>::iterator it = predB.begin(); it != predB.end(); ++
it)
            {
                Configuration* A = * it;
                Temp.insert(A);
            }
        }
    }
}
void TestSimulationOf2DPDA()
{
    vector<Configuration*> vecConfigs;
    Configuration c0;
    Configuration c1;
    Configuration c2;
    Configuration c3;
    Configuration c4;
    Configuration c5;
    Configuration c6;
    Configuration c7;

```

```
Configuration c8;
Configuration c9;
Configuration c10;
Configuration c11;
c0.iNumberOfMove = 1;
c1.iNumberOfMove = 2;
c2.iNumberOfMove = 2;
c3.iNumberOfMove = 3;
c4.iNumberOfMove = 4;
c5.iNumberOfMove = 3;
c6.iNumberOfMove = 3;
c7.iNumberOfMove = 2;
c8.iNumberOfMove = 3;
c9.iNumberOfMove = 3;
c10.iNumberOfMove = 2;
c11.iNumberOfMove = 1;
c0.strConfigName = "c0";
c1.strConfigName = "c1";
c2.strConfigName = "c2";
c3.strConfigName = "c3";
c4.strConfigName = "c4";
c5.strConfigName = "c5";
c6.strConfigName = "c6";
c7.strConfigName = "c7";
c8.strConfigName = "c8";
c9.strConfigName = "c9";
c10.strConfigName = "c10";
c11.strConfigName = "c11";

vecConfigs.push_back(&c0);
vecConfigs.push_back(&c1);
vecConfigs.push_back(&c2);
vecConfigs.push_back(&c3);
vecConfigs.push_back(&c4);
vecConfigs.push_back(&c5);
vecConfigs.push_back(&c6);
vecConfigs.push_back(&c7);
vecConfigs.push_back(&c8);
vecConfigs.push_back(&c9);
vecConfigs.push_back(&c10);
vecConfigs.push_back(&c11);
bool f = SimulationOf2DPDA(vecConfigs);
cout << "result: " << (f == true ? "true" : "false") << endl;
}
int main()
{
    TestSimulationOf2DPDA();
    system("pause");
    return 0;
}
```


参 考 文 献

- ADEL'SON-VEL'SKII, G. M., AND Y. M. LANDIS [1962]. "An algorithm for the organization of information," *Dokl. Akad. Nauk SSSR* 146, 263-266 (in Russian). English translation in *Soviet Math. Dokl.* 3 (1962), 1259-1262.
- AHO, A. V. (ed.) [1973]. *Currents in the Theory of Computing*, Prentice-Hall, Englewood Cliffs, N.J.
- AHO, A. V., M. R. GAREY, AND J. D. ULLMAN [1972]. "The transitive reduction of a directed graph," *SIAM J. Computing* 1:2, 131-137.
- AHO, A. V., D. S. HIRSCHBERG, AND J. D. ULLMAN [1974]. "Bounds on the complexity of the maximal common subsequence problem." Conference Record, IEEE 15th Annual Symposium on Switching and Automata Theory.
- AHO, A. V., J. E. HOPCROFT, AND J. D. ULLMAN [1968]. "Time and tape complexity of pushdown automaton languages," *Information and Control* 13:3, 186-206.
- AHO, A. V., J. E. HOPCROFT, AND J. D. ULLMAN [1974]. "On finding lowest common ancestors in trees," *SIAM J. Computing*, to appear.
- AHO, A. V., K. STEIGLITZ, AND J. D. ULLMAN [1974]. "Evaluating polynomials at fixed sets of points," Bell Laboratories, Murray Hill, N. J.
- AHO, A. V., AND J. D. ULLMAN [1972]. *The Theory of Parsing, Translation, and Compiling*, Volume 1: *Parsing*, Prentice-Hall, Englewood Cliffs, N.J.
- AHO, A. V., AND J. D. ULLMAN [1973]. *The Theory of Parsing, Translation, and Compiling*, Volume 2: *Compiling*, Prentice-Hall, Englewood Cliffs, N.J.
- ARLAZAROV, V. L., E. A. DINIC, M. A. KRONROD, AND I. A. FARADZEV [1970]. "On economical construction of the transitive closure of a directed graph," *Dokl. Akad. Nauk SSSR* 194, 487-488 (in Russian). English translation in *Soviet Math. Dokl.* 11:5, 1209-1210.
- BELAGA, E. C. [1961]. "On computing polynomials in one variable with initial preconditioning of the coefficients," *Problemi Kibernetiki* 5, 7-15.
- BELLMAN, R. E. [1957]. *Dynamic Programming*, Princeton University Press, Princeton, N.J.
- BERGE, C. [1958]. *The Theory of Graphs and its Applications*, Wiley, New York.
- BIRKHOFF, G., AND T. BARTEE [1970]. *Modern Applied Algebra*, McGraw-Hill, New York.
- BLUESTEIN, L. I. [1970]. "A linear filtering approach to the computation of the discrete Fourier transform," *IEEE Trans. on Electroacoustics AU-18:4*, 451-455. Also in Rabiner and Rader [1972], pp. 317-321.
- BLUM, M. [1967]. "A machine independent theory of the complexity of recursive functions," *J. ACM* 14:2, 322-336.
- BLUM, M., R. W. FLOYD, V. R. PRATT, R. L. RIVEST, AND R. E. TARJAN [1972]. "Time bounds for selection," *J. Computer and System Sciences* 7:4, 448-461.
- BOOK, R. V. [1972]. "On languages accepted in polynomial time," *SIAM J. Computing* 1:4, 281-287.
- BOOK, R. V. [1974]. "Comparing complexity classes," *J. Computer and System Sciences*, to appear.

- BORODIN, A. B. [1973a]. "Computational complexity—theory and practice," in Aho [1973], pp. 35–89.
- BORODIN, A. B. [1973b]. "On the number of arithmetics required to compute certain functions—circa May 1973," in Traub [1973], pp. 149–180.
- BORODIN, A. B., AND S. A. COOK [1974]. "On the number of additions to compute specific polynomials," *Proc. 6th Annual ACM Symposium on Theory of Computing*, 342–347.
- BORODIN, A. B., AND I. MUNRO [1971]. "Evaluating polynomials at many points," *Information Processing Letters* 1:2, 66–68.
- BORODIN, A. B., AND I. MUNRO [1975]. *The Computational Complexity of Algebraic and Numeric Problems*, to appear, American Elsevier, N.Y.
- BROWN, W. S. [1971]. "On Euclid's algorithm and greatest common divisors," *J. ACM* 18:4, 478–504.
- BROWN, W. S. [1973]. *ALTRAN User's Manual* (3rd edition), Bell Laboratories, Murray Hill, N.J.
- BRUNO, J. L., AND R. SETHI [1974]. "Register allocation for a one-register machine," *Proc. 8th Annual Princeton Conference on Information Sciences and Systems*.
- BUNCH, J., AND J. E. HOPCROFT [1974]. "Triangular factorization and inversion by fast matrix multiplication," *Math. Comp.* 28:125, 231–236.
- BURKHARD, W. A. [1973]. "Nonrecursive tree traversal algorithms," *Proc. 7th Annual Princeton Conference on Information Sciences and Systems*, 403–405.
- COLLINS, G. E. [1973]. "Computer algebra of polynomials and rational functions," *American Math. Monthly* 80:7, 725–754.
- CONSTABLE, R. L., H. B. HUNT III, AND S. K. SAHNI [1974]. "On the computational complexity of scheme equivalence," *Proc. 8th Annual Princeton Conference on Information Sciences and Systems*.
- COOK, S. A. [1971a]. "Linear time simulation of deterministic two-way pushdown automata," *Proc. IFIP Congress 71*, TA-2. North-Holland, Amsterdam, 172–179.
- COOK, S. A. [1971b]. "The complexity of theorem proving procedures," *Proc. 3rd Annual ACM Symposium on Theory of Computing*, 151–158.
- COOK, S. A. [1973]. "A hierarchy for nondeterministic time complexity," *J. Computer and System Sciences* 7:4, 343–353.
- COOK, S. A., AND S. O. AANDERAA [1969]. "On the minimum complexity of functions," *Trans. Amer. Math. Soc.* 142, 291–314.
- COOK, S. A., AND R. A. RECKHOW [1973]. "Time-bounded random access machines," *J. Computer and System Sciences* 7:4, 354–375.
- COOK, S. A., AND R. A. RECKHOW [1974]. "On the lengths of proofs in the propositional calculus," *Proc. 6th Annual ACM Symposium on Theory of Computing*, 135–148.
- COOLEY, J. M., P. A. LEWIS, AND P. D. WELCH [1967]. "History of the fast Fourier transform," *Proc. IEEE* 55, 1675–1677.

- COOLEY, J. M., AND J. W. TUKEY [1965]. "An algorithm for the machine calculation of complex Fourier series," *Math. Comp.* 19, 297-301.
- CRANE, C. A. [1972]. "Linear lists and priority queues as balanced binary trees," Ph.D. Thesis, Stanford University.
- DANIELSON, G. C., AND C. LANCZOS [1942]. "Some improvements in practical Fourier analysis and their application to X-ray scattering from liquids," *J. Franklin Institute* 233, 365-380 and 435-452.
- DANTZIG, G. B., W. O. BLATTNER, AND M. R. RAO [1966]. "All shortest routes from a fixed origin in a graph," *Theory of Graphs*, International Symposium, Rome, July 1966. Proceedings published by Gordon and Breach, New York, 1967, pp. 85-90.
- DAVIS, M. [1958]. *Computability and Unsolvability*, McGraw-Hill, New York.
- DIJKSTRA, E. W. [1959]. "A note on two problems in connexion with graphs," *Numerische Mathematik* 1, 269-271.
- DIVETTI, L., AND A. GRASSELLI [1968]. "On the determination of minimum feedback arc and vertex sets," *IEEE Trans. Circuit Theory* CT-15:1, 86-89.
- EHRENFEUCHT, A., AND H. P. ZEIGER [1974]. "Complexity measures for regular expressions," *Proc. 6th Annual ACM Symposium on Theory of Computing*, 75-79.
- ELGOT, C. C., AND A. ROBINSON [1964]. "Random access stored program machines," *J. ACM* 11:4, 365-399.
- EVE, J. [1964]. "The evaluation of polynomials," *Numerische Mathematik* 6, 17-21.
- EVEN, S. [1973]. "An algorithm for determining whether the connectivity of a graph is at least k ," TR-73-184, Dept. of Computer Science, Cornell University, Ithaca, N.Y.
- FIDUCCIA, C. M. [1971]. "Fast matrix multiplication," *Proc. 3rd Annual ACM Symposium on Theory of Computing*, 45-49.
- FIDUCCIA, C. M. [1972]. "Polynomial evaluation via the division algorithm—the fast Fourier transform revisited," *Proc. 4th Annual ACM Symposium on Theory of Computing*, 88-93.
- FISCHER, M. J. [1972]. "Efficiency of equivalence algorithms," in Miller and Thatcher [1972], pp. 153-168.
- FISCHER, M. J., AND A. R. MEYER [1971]. "Boolean matrix multiplication and transitive closure," *Conference Record, IEEE 12th Annual Symposium on Switching and Automata Theory*, 129-131.
- FISCHER, M. J., AND M. S. PATERSON [1974]. "String-matching and other products." Project MAC Technical Memorandum 41, MIT, Cambridge, Mass.
- FISCHER, M. J., AND M. O. RABIN [1974]. "The complexity of theorem proving procedures," Project MAC Report, MIT, Cambridge, Mass.
- FLOYD, R. W. [1962]. "Algorithm 97: shortest path," *Comm. ACM* 5:6, 345.
- FLOYD, R. W. [1964]. "Algorithm 245: treesort 3," *Comm. ACM* 7:12, 701.
- FLOYD, R. W., AND R. L. RIVEST [1973]. "Expected time bounds for selection," Computer Science Dept., Stanford University.

- FORD, L. R., AND S. M. JOHNSON [1959]. "A tournament problem," *Amer. Math. Monthly* 66, 387-389.
- FRAZER, W. D., AND A. C. MCKELLAR [1970]. "Samplesort: a sampling approach to minimal storage tree sorting," *J. ACM* 17:3, 496-507.
- FURMAN, M. E. [1970]. "Application of a method of fast multiplication of matrices in the problem of finding the transitive closure of a graph," *Dokl. Akad. Nauk SSSR* 194, 524 (in Russian). English translation in *Soviet Math. Dokl.* 11:5, 1252.
- GALE, D., AND R. M. KARP [1970]. "A phenomenon in the theory of sorting," *Conference Record, IEEE 11th Annual Symposium on Switching and Automata Theory*, 51-59.
- GALLER, B. A., AND M. J. FISCHER [1964]. "An improved equivalence algorithm," *Comm. ACM* 7:5, 301-303.
- GAREY, M. R., D. S. JOHNSON, AND L. STOCKMEYER [1974]. "Some simplified polynomial complete problems," *Proc. 6th Annual ACM Symposium on Theory of Computing*, 47-63.
- GENTLEMAN, W. M. [1972]. "Optimal multiplication chains for computing powers of polynomials," *Math. Comp.* 26:120, 935-940.
- GENTLEMAN, W. M., AND S. C. JOHNSON [1973]. "Analysis of algorithms, a case study: determinants of polynomials," *Proc. 5th Annual ACM Symposium on Theory of Computing*, 135-141.
- GENTLEMAN, W. M., AND G. SANDE [1966]. "Fast Fourier transforms for fun and profit," *Proc. AFIPS 1966 Fall Joint Computer Conference*, v. 29, Spartan, Washington, D.C., pp. 563-578.
- GILBERT, E. N., AND E. F. MOORE [1959]. "Variable length encodings," *Bell System Technical J.* 38:4, 933-968.
- GODBOLE, S. S. [1973]. "On efficient computation of matrix chain products," *IEEE Transactions on Computers* C-22:9, 864-866.
- GOOD, I. J. [1958]. "The interaction algorithm and practical Fourier series," *J. Royal Statistics Soc., Ser. B*, 20, 361-372. Addendum 22 (1960), 372-375.
- GRAHAM, R. L. [1972]. "An efficient algorithm for determining the convex hull of a planar set," *Information Processing Letters* 1:4, 132-133.
- GRAY, J. N., M. A. HARRISON, AND O. H. IBARRA [1967]. "Two way pushdown automata," *Information and Control* 11:3, 30-70.
- HADIAN, A., AND M. SOBEL [1969]. "Selecting the n th largest using binary errorless comparisons," *Tech. Rept. 121*, Department of Statistics, University of Minnesota, Minneapolis.
- HALL, A. D. [1971]. "The ALTRAN system for rational manipulation—a survey," *Comm. ACM* 14:8, 517-521.
- HARARY, F. [1969]. *Graph Theory*, Addison-Wesley, Reading, Mass.
- HARPER, L. H., AND J. E. SAVAGE [1972]. "On the complexity of the marriage problem," in *Advances in Mathematics* 9:3, 299-312.
- HARTMANIS, J. [1971]. "Computational complexity of random access stored program machines," *Mathematical Systems Theory* 5:3, 232-245.

- HARTMANIS, J., AND J. E. HOPCROFT [1971]. "An overview of the theory of computational complexity," *J. ACM* 18:3, 444-475.
- HARTMANIS, J., P. M. LEWIS II, AND R. E. STEARNS [1965]. "Classification of computations by time and memory requirements," *Proc. IFIP Congress 65*, Spartan, N.Y., 31-35.
- HARTMANIS, J., AND R. E. STEARNS [1965]. "On the computational complexity of algorithms," *Trans. Amer. Math. Soc.* 117, 285-306.
- HECHT, M. S. [1973]. "Global data flow analysis of computer programs," Ph.D. Thesis, Dept. of Electrical Engineering, Princeton University.
- HEINDEL, L. E., AND E. HOROWITZ [1971]. "On decreasing the computing time for modular arithmetic," *Conference Record, IEEE 12th Annual Symposium on Switching and Automata Theory*, 126-128.
- HENNIE, F. C., AND R. E. STEARNS [1966]. "Two tape simulation of multitape machines," *J. ACM* 13:4, 533-546.
- HIRSCHBERG, D. S. [1973]. "A linear space algorithm for computing maximal common subsequences," TR-138, Computer Science Laboratory, Dept. of Electrical Engineering, Princeton University, Princeton, N.J.
- HOARE, C. A. R. [1962]. "Quicksort," *Computer J.* 5:1 10-15.
- HOHN, F. E. [1958]. *Elementary Matrix Algebra*, Macmillan, New York.
- HOPCROFT, J. E. [1971]. "An $n \log n$ algorithm for minimizing states in a finite automaton," in Kohavi and Paz [1971], pp. 189-196.
- HOPCROFT, J. E., AND R. M. KARP [1971]. "An algorithm for testing the equivalence of finite automata," TR-71-114, Dept. of Computer Science, Cornell University, Ithaca, N.Y.
- HOPCROFT, J. E., AND L. R. KERR [1971]. "On minimizing the number of multiplications necessary for matrix multiplication," *SIAM J. Applied Math.* 20:1, 30-36.
- HOPCROFT, J. E., AND J. MUSINSKI [1973]. "Duality in determining the complexity of noncommutative matrix multiplication," *Proc. 5th Annual ACM Symposium on Theory of Computing*, 73-87.
- HOPCROFT, J. E., AND R. E. TARJAN [1973a]. "Efficient planarity testing," TR 73-165, Dept. of Computer Science, Cornell University, Ithaca, N.Y., also *J. ACM*, to appear.
- HOPCROFT, J. E., AND R. E. TARJAN [1973b]. "Dividing a graph into triconnected components," *SIAM J. Computing* 2:3.
- HOPCROFT, J. E., AND R. E. TARJAN [1973c]. "Efficient algorithms for graph manipulation," *Comm. ACM* 16:6, 372-378.
- HOPCROFT, J. E., AND J. D. ULLMAN [1969]. *Formal Languages and Their Relation to Automata*, Addison-Wesley, Reading, Mass.
- HOPCROFT, J. E., AND J. D. ULLMAN [1973]. "Set merging algorithms," *SIAM J. Computing*, 2:4, 294-303.
- HOPCROFT, J. E., AND J. K. WONG [1974]. "A linear time algorithm for isomorphism of planar graphs," *Proc. 6th Annual ACM Symposium on Theory of Computing*, 172-184.

- HOROWITZ, E. [1972]. "A fast method for interpolation using preconditioning," *Information Processing Letters* 1:4, 157-163.
- HORVATH, E. C. [1974]. "Some efficient stable sorting algorithms," *Proc. 6th Annual ACM Symposium on Theory of Computing*, 194-215.
- HU, T. C. [1968]. "A decomposition algorithm for shortest paths in a network," *Operations Res.* 16, 91-102.
- HU, T. C., AND A. C. TUCKER [1971]. "Optimum binary search trees," *SIAM J. Applied Math.* 21:4, 514-532.
- HUNT, H. B., III [1973a]. "On the time and tape complexity of languages," *Proc. 5th Annual ACM Symposium on Theory of Computing*, 10-19. See also TR-73-182, Dept. of Computer Science, Cornell University, Ithaca, N.Y.
- HUNT, H. B., III [1973b]. "The equivalence problem for regular expressions with intersection is not polynomial in tape," TR 73-156, Dept. of Computer Science, Cornell University, Ithaca, N.Y.
- HUNT, H. B., III [1974]. "Stack languages, Rice's theorem and a natural exponential complexity gap," TR-142, Computer Science Laboratory, Dept. of Electrical Engineering, Princeton University, Princeton, N.J.
- HUNT, H. B., III, AND D. J. ROSENKRANTZ [1974]. "Computational parallels between the regular and context-free languages," *Proc. 6th Annual ACM Symposium on Theory of Computing*, 64-74.
- IBARRA, O. H. [1972]. "A note concerning nondeterministic tape complexities," *J. ACM* 19:4, 608-612.
- IBARRA, O. H., AND S. K. SAHNI [1973]. "Polynomial complete fault detection problems," Technical Report 74-3, Computer, Information and Control Sciences, University of Minnesota, Minneapolis.
- JONES, N. D. [1973]. "Reducibility among combinatorial problems in $\log n$ space," *Proc. 7th Annual Princeton Conference on Information Sciences and Systems*, 547-551.
- JOHNSON, D. B. [1973]. "Algorithms for shortest paths," Ph.D. Thesis, Dept. of Computer Science, Cornell University, Ithaca, New York.
- JOHNSON, S. C. [1974]. "Sparse polynomial arithmetic," Bell Laboratories, Murray Hill, N.J.
- KARATSUBA, A., AND Y. OFMAN [1962]. "Multiplication of multidigit numbers on automata," *Dokl. Akad. Nauk SSSR* 145, 293-294 (in Russian).
- KASAMI, T. [1965]. "An efficient recognition and syntax algorithm for context-free languages," *Scientific Report AFCRL-65-758*, Air Force Cambridge Research Laboratory, Bedford, Mass.
- KARP, R. M. [1972]. "Reducibility among combinatorial problems," in Miller and Thatcher [1972], pp. 85-104.
- KARP, R. M., R. E. MILLER, AND A. L. ROSENBERG [1972]. "Rapid identification of repeated patterns in strings, trees, and arrays," *Proc. 4th Annual ACM Symposium on Theory of Computing*, 125-136.

- KEDEM, Z. [1974]. "On the number of multiplications and divisions required to compute certain rational functions," *Proc. 6th Annual ACM Symposium on Theory of Computing*, 334-341.
- KERR, L. R. [1970]. "The effect of algebraic structure on the computational complexity of matrix multiplications," Ph.D. Thesis, Cornell University, Ithaca, N.Y.
- KIRKPATRICK, D. [1972]. "On the additions necessary to compute certain functions," *Proc. 4th Annual ACM Symposium on Theory of Computing*, 94-101.
- KIRKPATRICK, D. [1974]. "Determining graph properties from matrix representations," *Proc. 6th Annual ACM Symposium on Theory of Computing*, 84-90.
- KISLITSYN, S. S. [1964]. "On the selection of the k th element of an ordered set by pairwise comparison," *Sibirsk. Mat. Zh.* 5, 557-564.
- KLEENE, S. C. [1956]. "Representation of events in nerve nets and finite automata," in *Automata Studies* (Shannon and McCarthy, eds.), Princeton University Press, pp. 3-40.
- KNUTH, D. E. [1968]. *Fundamental Algorithms*, Addison-Wesley, Reading, Mass.
- KNUTH, D. E. [1969]. *Seminumerical Algorithms*, Addison-Wesley, Reading, Mass.
- KNUTH, D. E. [1971]. "Optimum binary search trees," *Acta Informatica* 1, 14-25.
- KNUTH, D. E. [1973a]. *Sorting and Searching*, Addison-Wesley, Reading, Mass.
- KNUTH, D. E. [1973b]. "Notes on pattern matching," University of Trondheim, Norway.
- KNUTH, D. E., AND V. R. PRATT [1971]. "Automata theory can be useful," Stanford University; Stanford, California.
- KOHAVER, Z., AND A. PAZ (eds.) [1971]. *Theory of Machines and Computations*, Academic Press, New York.
- KRUSKAL, J. B., JR. [1956]. "On the shortest spanning subtree of a graph and the traveling salesman problem," *Proc. Amer. Math. Soc.* 7:1, 48-50.
- KUNG, H. T. [1973]. "Fast evaluation and interpolation," Dept. of Computer Science, Carnegie-Mellon University, Pittsburgh.
- LEWIS, P. M., II, R. E. STEARNS, AND J. HARTMANIS [1965]. "Memory bounds for recognition of context-free and context-sensitive languages," *Conference Record, IEEE 6th Annual Symposium on Switching Circuit Theory and Logic Design*, 191-202.
- LIPSON, J. [1971]. "Chinese remainder and interpolation algorithms," *Proc. 2nd Symposium on Symbolic and Algebraic Manipulation*, 372-391.
- LIU, C. L. [1968]. *Introduction to Combinatorial Mathematics*, McGraw-Hill, New York.
- LIU, C. L. [1972]. "Analysis and synthesis of sorting algorithms," *SIAM J. Computing* 1:4, 290-304.
- MACLANE, S., AND G. BIRKHOFF [1967]. *Algebra*, Macmillan, New York.
- MCCAUGHTON, R., AND H. YAMADA [1960]. "Regular expressions and state graphs for automata," *IRE Trans. on Electronic Computers* 9:1, 39-47. Reprinted in Moore [1964], pp. 157-174.

- MEYER, A. R. [1972]. "Weak monadic second order theory of successor is not elementary recursive," Project MAC Report, MIT, Cambridge, Mass.
- MEYER, A. R., AND L. STOCKMEYER [1972]. "The equivalence problem for regular expressions with squaring requires exponential space," *Conference Record, IEEE 13th Annual Symposium on Switching and Automata Theory*, 125-129.
- MEYER, A. R., AND L. STOCKMEYER [1973]. "Nonelementary word problems in automata and logic," *Proc. AMS Symposium on Complexity of Computation*, April 1973.
- MILLER, R. E., AND J. W. THATCHER (eds.) [1972]. *Complexity of Computer Computations*, Plenum Press, New York.
- MINSKY, M. [1967]. *Computation: Finite and Infinite Machines*, Prentice-Hall, Englewood Cliffs, N.J.
- MOENCK, R. [1973]. "Fast Computation of GCD's," *Proc. 5th Annual ACM Symposium on Theory of Computing*, 142-151.
- MOENCK, R., AND A. B. BORODIN [1972]. "Fast modular transforms via division," *Conference Record, IEEE 13th Annual Symposium on Switching and Automata Theory*, 90-96.
- MOORE, E. F. (ed.) [1964]. *Sequential Machines: Selected Papers*, Addison-Wesley, Reading, Mass.
- MORGENSTERN, J. [1973]. "Note on a lower bound of the linear complexity of the fast Fourier transform," *J. ACM* 20:2, 305-306.
- MORRIS, R. [1968]. "Scatter storage techniques," *Comm. ACM* 11:1, 35-44.
- MORRIS, J. H., JR., AND V. R. PRATT [1970]. "A linear pattern matching algorithm," Technical Report No. 40, Computing Center, University of California, Berkeley.
- MOTZKIN, T. S. [1955]. "Evaluation of polynomials and evaluation of rational functions," *Bull. Amer. Math. Soc.* 61, 163.
- MUNRO, I. [1971]. "Efficient determination of the transitive closure of a directed graph," *Information Processing Letters* 1:2, 56-58.
- MURAOKA, Y., AND D. J. KUCK [1973]. "On the time required for a sequence of matrix products," *Comm. ACM* 16:1, 22-26.
- NECIPORUK, E. I. [1966]. "A Boolean function," *Dokl. Akad. Nauk SSSR* 169:4 (in Russian). English translation in *Soviet Math. Dokl.* 7 (1966), 999-1000.
- NICHOLSON, P. J. [1971]. "Algebraic theory of finite Fourier transforms," *J. Computer and System Sciences* 5:5, 524-547.
- NIEVERGELT, J., AND E. M. REINGOLD [1973]. "Binary search trees of bounded balance," *SIAM J. Computing* 2:1, 33-43.
- OFMAN, Y. [1962]. "On the algorithmic complexity of discrete functions," *Dokl. Akad. Nauk SSSR* 145, 48-51 (in Russian). English translation in *Soviet Physics Dokl.* 7 (1963), 589-591.
- OSTROWSKI, A. M. [1954]. "On two problems in abstract algebra connected with Horner's rule," *Studies Presented to R. von Mises*, Academic Press, N.Y.
- PAN, V. Y. [1966]. "Methods of computing values of polynomials," *Russian Mathematical Surveys* 21:1, 105-136.

- PATERSON, M. S., M. J. FISCHER, AND A. R. MEYER [1974]. "An improved overlap argument for on-line multiplication," Project MAC Technical Memorandum 40, MIT, Cambridge, Mass.
- POHL, I. [1972]. "A sorting problem and its complexity," *Comm. ACM* 15:6, 462-464.
- PRATT, T. W. [1975]. *Programming Language Design and Implementation*, to appear, Prentice-Hall, Englewood Cliffs, N.J.
- PRATT, V. R. [1974]. "The power of negative thinking in multiplying Boolean matrices," *Proc. 6th Annual ACM Symposium on Theory of Computing*, 80-83.
- PRATT, V. R., AND F. F. YAO [1973]. "On lower bounds for computing the i th largest element," *Conference Record, IEEE 14th Annual Symposium on Switching and Automata Theory*, 70-81.
- PRIM, R. C. [1957]. "Shortest connection networks and some generalizations," *Bell System Technical J.*, 1389-1401.
- RABIN, M. O. [1963]. "Real-time computation," *Israel J. Math.* 1, 203-211.
- RABIN, M. O. [1972]. "Proving simultaneous positivity of linear forms," *J. Computer and System Sciences* 6:6, 639-650.
- RABIN, M. O., AND D. SCOTT [1959]. "Finite automata and their decision problems," *IBM J. Research and Development* 3, 114-125. Reprinted in Moore [1964], pp. 63-91.
- RABIN, M. O., AND S. WINOGRAD [1971]. "Fast evaluation of polynomials by rational preparation," *IBM Technical Report RC 3645*, Yorktown Heights, N.Y.
- RABINER, L. R., AND C. M. RADER (eds.) [1972]. *Digital Signal Processing*, IEEE Press, New York.
- RABINER, L. R., R. W. SCHAFER, AND C. M. RADER [1969]. "The chirp z -transform and its applications," *Bell System Technical J.* 48:3, 1249-1292. Reprinted in Rabiner and Rader [1972], pp. 322-328.
- RADER, C. M. [1968]. "Discrete Fourier transform when the number of data points is prime," *Proc. IEEE* 56, 1107-1108.
- REINGOLD, E. M. [1972]. "On the optimality of some set algorithms," *J. ACM* 19:4, 649-659.
- ROGERS, H., JR. [1967]. *Theory of Recursive Functions and Effective Computability*, McGraw-Hill, New York.
- ROUNDS, W. C. [1973]. "Complexity of recognition in intermediate level languages," *Conference Record, IEEE 14th Annual Symposium on Switching and Automata Theory*, 145-158.
- RUNGE, C., AND H. KÖNIG [1924]. *Die Grundlehren der mathematischen Wissenschaften*, v. 11, Springer, Berlin.
- SAHNI, S. K. [1972]. "Some related problems from network flows, game theory and integer programming," *Conference Record, IEEE 13th Annual Symposium on Switching and Automata Theory*, 130-138.
- SAVITCH, W. J. [1970]. "Relationship between nondeterministic and deterministic tape complexities," *J. Computer and System Sciences* 4:2, 177-192.

- SAVITCH, W. J. [1971]. "Maze recognizing automata." *Proc. 4th Annual ACM Symposium on Theory of Computing*, 151-156.
- SCHÖNHAGE, A. [1971]. "Schnelle Berechnung von Kettenbruchentwicklungen," *Acta Informatica* 1, 139-144.
- SCHÖNHAGE, A., AND V. STRASSEN [1971]. "Schnelle Multiplikation grosser Zahlen," *Computing* 7, 281-292.
- SEIFERAS, J. J., M. J. FISCHER, AND A. R. MEYER [1973]. "Refinements of nondeterministic time and space hierarchies," *Conference Record, IEEE 14th Annual Symposium on Switching and Automata Theory*, 130-137.
- SETHI, R. [1973]. "Complete register allocation problems," *Proc. 5th Annual ACM Symposium on Theory of Computing*, 182-195.
- SHEPHERDSON, J. C., AND H. E. STURGIS [1963]. "Computability of recursive functions," *J. ACM* 10:2, 217-255.
- SIEVEKING, M. [1972]. "An algorithm for division of power series," *Computing* 10, 153-156.
- SINGLETON, R. C. [1969]. "Algorithm 347: An algorithm for sorting with minimal storage." *Comm. ACM* 12:3, 185-187.
- SLOANE, N. J. A. [1973]. *A Handbook of Integer Sequences*, Academic Press, N.Y.
- SPIRA, P. M. [1973]. "A new algorithm for finding all shortest paths in a graph of positive arcs in average time $O(n^2 \log^2 n)$," *SIAM J. Computing* 2:1, 28-32.
- SPIRA, P. M., AND A. PAN [1973]. "On finding and updating shortest paths and spanning trees," *Conference Record, IEEE 14th Annual Symposium on Switching and Automata Theory*, 82-84.
- STEARNS, R. E., AND D. J. ROSENKRANTZ [1969]. "Table machine simulation," *Conference Record, IEEE 10th Annual Symposium on Switching and Automata Theory*, 118-128.
- STOCKMEYER, L. [1973]. "Planar 3-colorability is polynomial complete," *SIGACT News* 5:3, 19-25.
- STOCKMEYER, L., AND A. R. MEYER [1973]. "Word problems requiring exponential time," *Proc. 5th Annual ACM Symposium on Theory of Computing*, 1-9.
- STONE, H. S. [1972]. *Introduction to Computer Organization and Data Structures*, McGraw-Hill, New York.
- STRASSEN, V. [1969]. "Gaussian elimination is not optimal," *Numerische Mathematik* 13, 354-356.
- STRASSEN, V. [1974]. "Schwer berechenbare Polynome mit rationalen Koeffizienten," *SIAM J. Computing*.
- TARJAN, R. E. [1972]. "Depth first search and linear graph algorithms," *SIAM J. Computing* 1:2, 146-160.
- TARJAN, R. E. [1973a]. "Finding dominators in directed graphs," *Proc. 7th Annual Princeton Conference on Information Sciences and Systems*, 414-418.
- TARJAN, R. E. [1973b]. "Testing flow graph reducibility," *Proc. 5th Annual ACM Symposium on Theory of Computing*, 96-107.

- TARJAN, R. E. [1974]. "On the efficiency of a good but not linear set merging algorithm," *J. ACM*, to appear.
- THOMPSON, K. [1968]. "Regular expression search algorithm," *Comm. ACM* 11:6, 419-422.
- TRAUB, J. F. (ed.) [1973]. *Complexity of Sequential and Parallel Numerical Algorithms*, Academic Press, New York.
- TURING, A. M. [1936]. "On computable numbers, with an application to the Entscheidungsproblem," *Proc. London Mathematical Soc. Ser. 2*, 42, 230-265. Corrections, *ibid.*, 43 (1937), 544-546.
- ULLMAN, J. D. [1973]. "Polynomial complete scheduling problems," *Proc. 4th Symposium on Operating System Principles*, 96-101.
- ULLMAN, J. D. [1974]. "Fast algorithms for the elimination of common subexpressions," *Acta Informatica*, 2:3, 191-213.
- VALIANT, L. G. [1974]. "General context-free recognition in less than cubic time." Dept. of Computer Science, Carnegie-Mellon University, Pittsburg, Pa.
- VARI, T. M. "Some complexity results for a class of Toeplitz matrices," unpublished memorandum, York Univ., Toronto, Ontario, Canada.
- WAGNER, R. A. [1974]. "A shortest path algorithm for edge-sparse graphs." Dept. of Systems and Information Science, Vanderbilt University, Nashville, Tennessee.
- WAGNER, R. A., AND M. J. FISCHER [1974]. "The string-to-string correction problem," *J. ACM* 21:1, 168-173.
- WARSHALL, S. [1962]. "A theorem on Boolean matrices," *J. ACM* 9:1, 11-12.
- WEINER, P. [1973]. "Linear pattern matching algorithms," *Conference Record, IEEE 14th Annual Symposium on Switching and Automata Theory*, 1-11.
- WILLIAMS, J. W. J. [1964]. "Algorithm 232: Heapsort," *Comm. ACM* 7:6, 347-348.
- WINOGRAD, S. [1965]. "On the time required to perform addition," *J. ACM* 12:2, 277-285.
- WINOGRAD, S. [1967]. "On the time required to perform multiplication," *J. ACM* 14:4, 793-802.
- WINOGRAD, S. [1970a]. "On the number of multiplications necessary to compute certain functions," *Comm. Pure and Applied Math.* 23, 165-179.
- WINOGRAD, S. [1970b]. "On the multiplication of 2×2 matrices," *IBM Research Report RC267*, January 1970.
- WINOGRAD, S. [1970c]. "The algebraic complexity of functions," *Proc. International Congress of Mathematicians* (1970), vol. 3, 283-288.
- WINOGRAD, S. [1973]. "Some remarks on fast multiplication of polynomials," in Traub [1973], pp. 181-196.
- YOUNGER, D. H. [1967]. "Recognition of context-free languages in time n^3 ," *Information and Control* 10:2, 189-208.